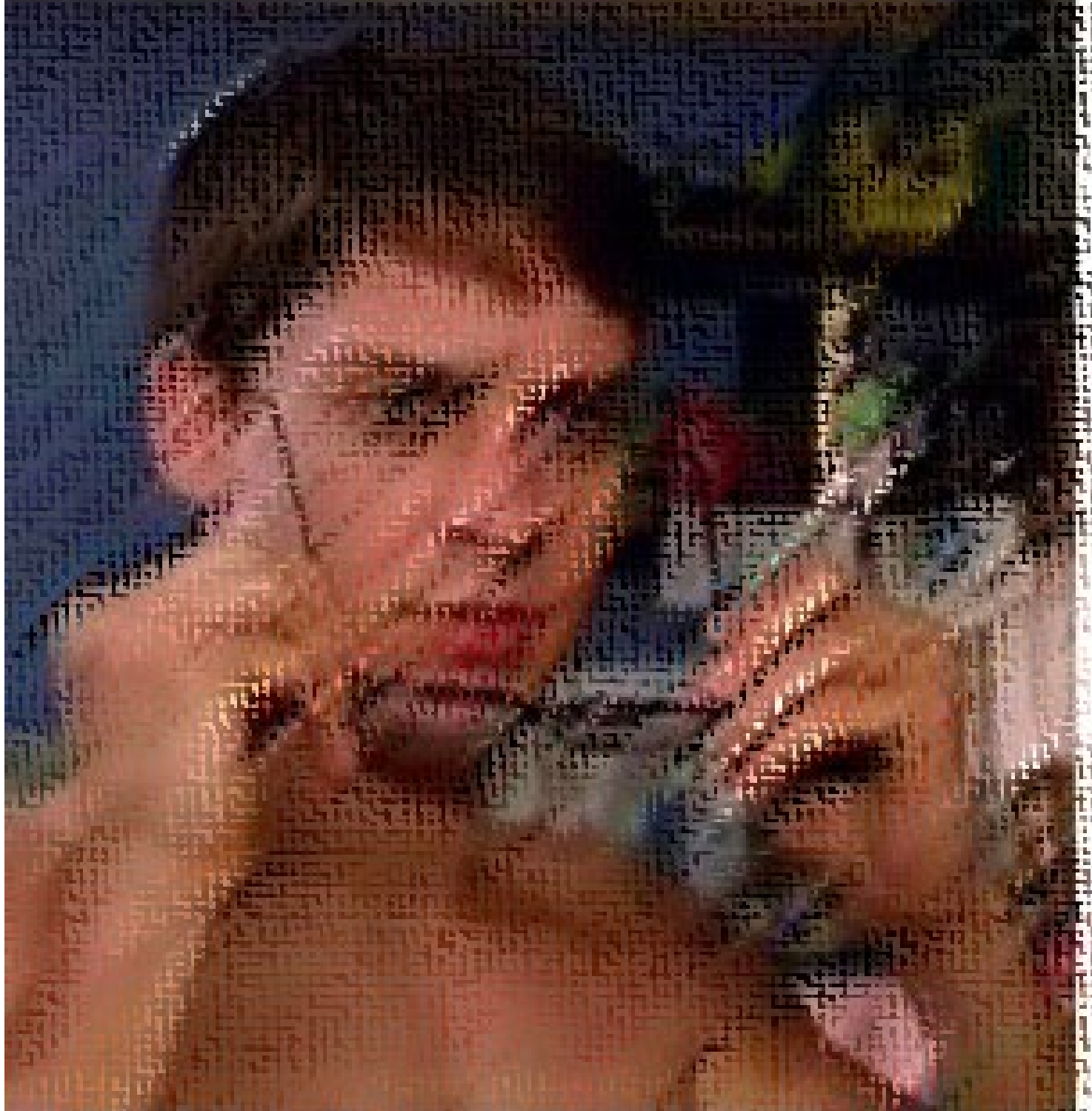


Lecture 2: Neural Networks



Some slides from Suren Jayasuriya (ASU), Phillip Isola (MIT)



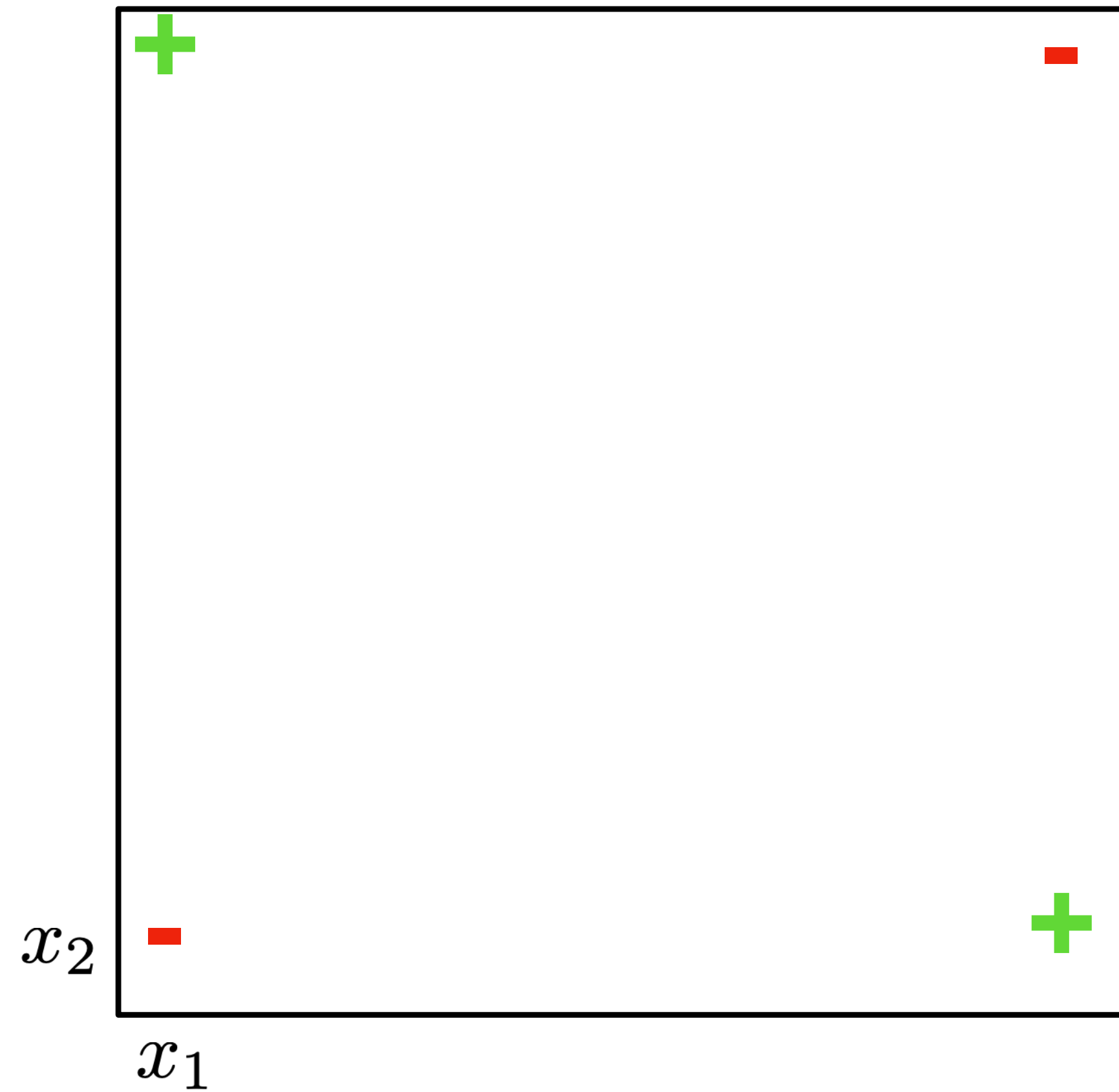


Artificial
Intelligence



$$\hat{y} = w^T x + b$$

Limitations to linear classifiers

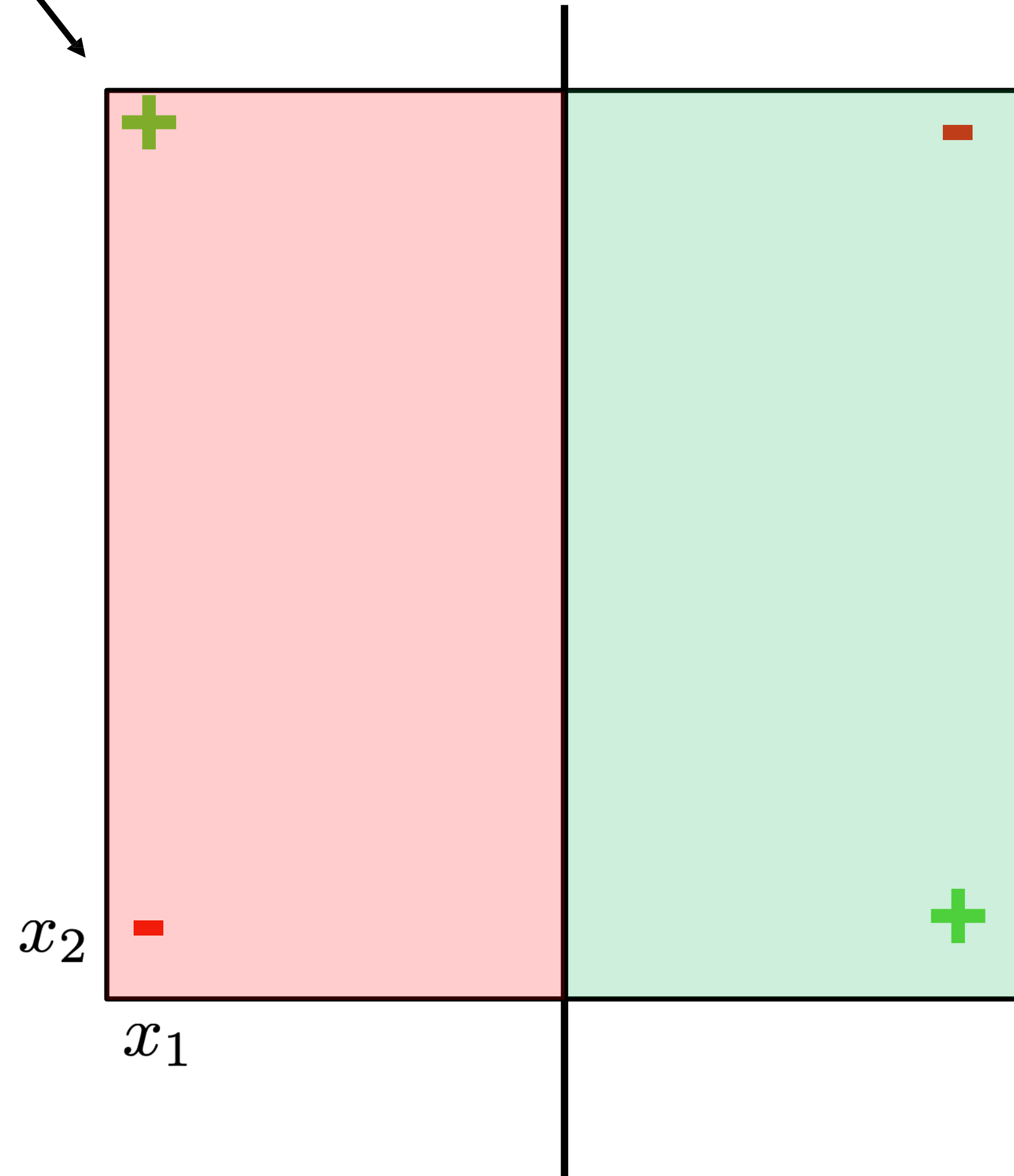


		x_2	
		0	1
x_1	0	0	1
	1	1	0
		XOR	

Limitations to linear classifiers

Wrong!

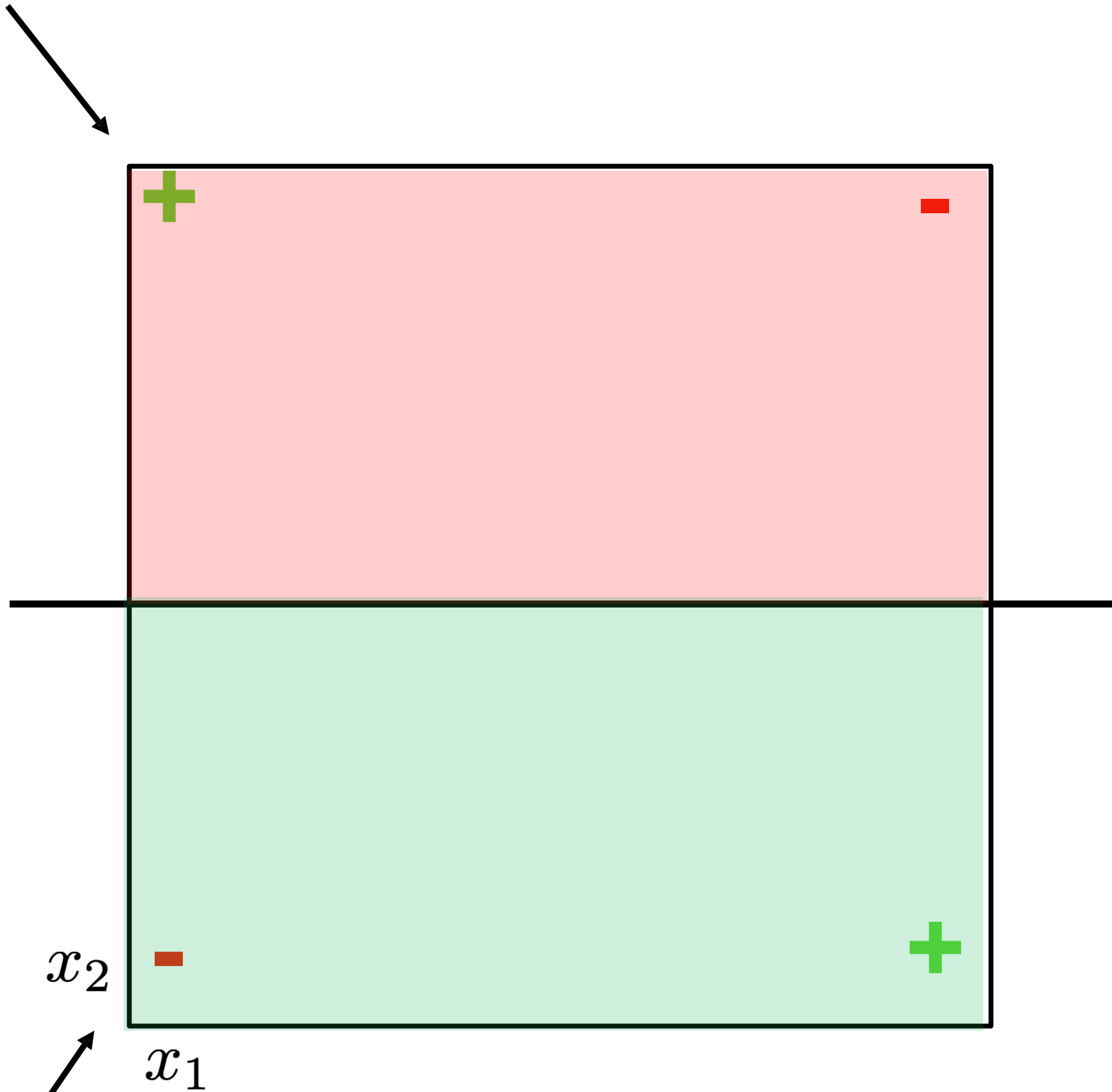
Wrong!



		x_2	
		0	1
x_1	0	0	1
	1	1	0
		XOR	

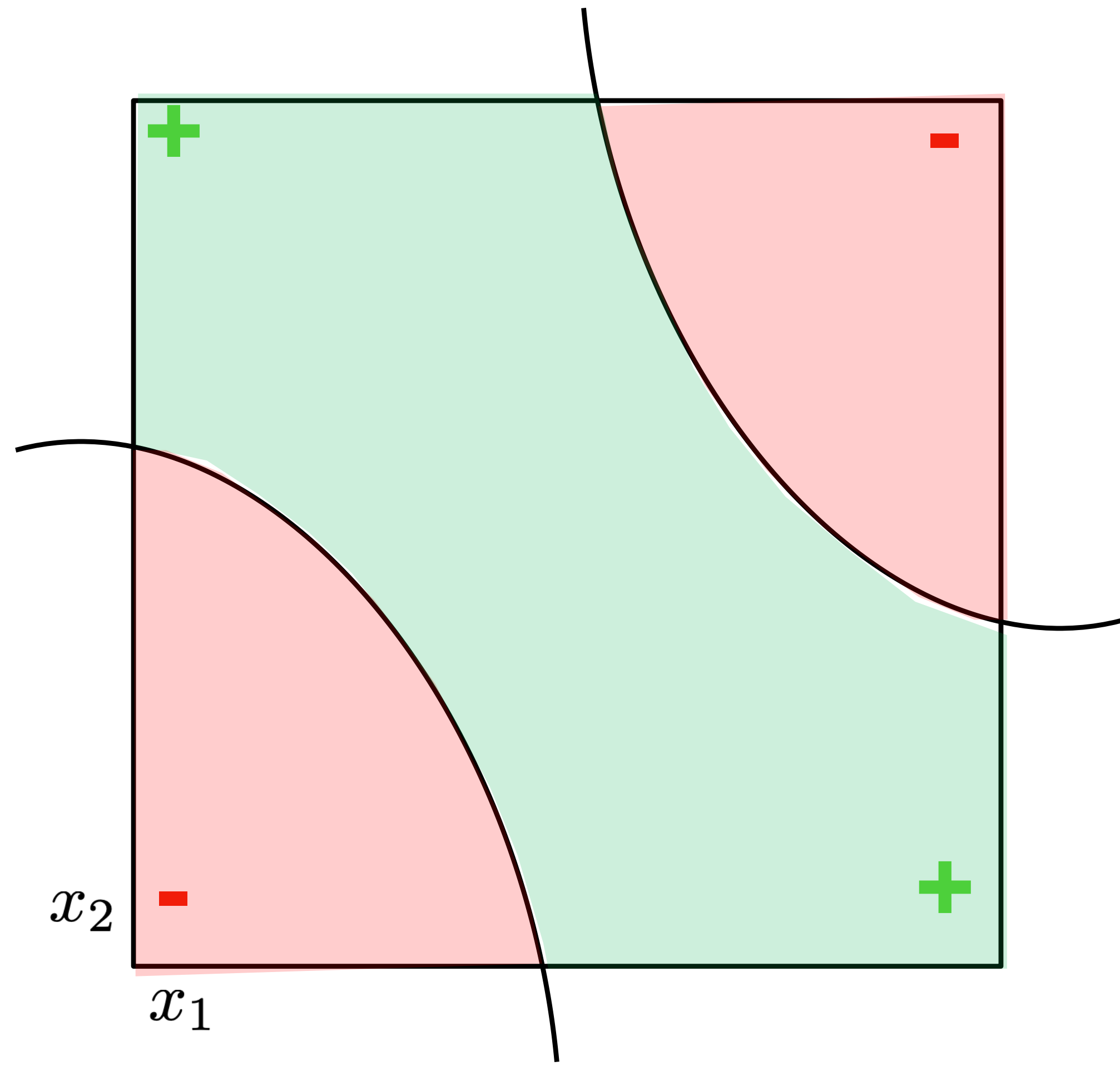
Limitations to linear classifiers

Wrong!



		x_2	
		0	1
x_1	0	0	1
	1	1	0
		XOR	

Goal: Non-linear decision boundary



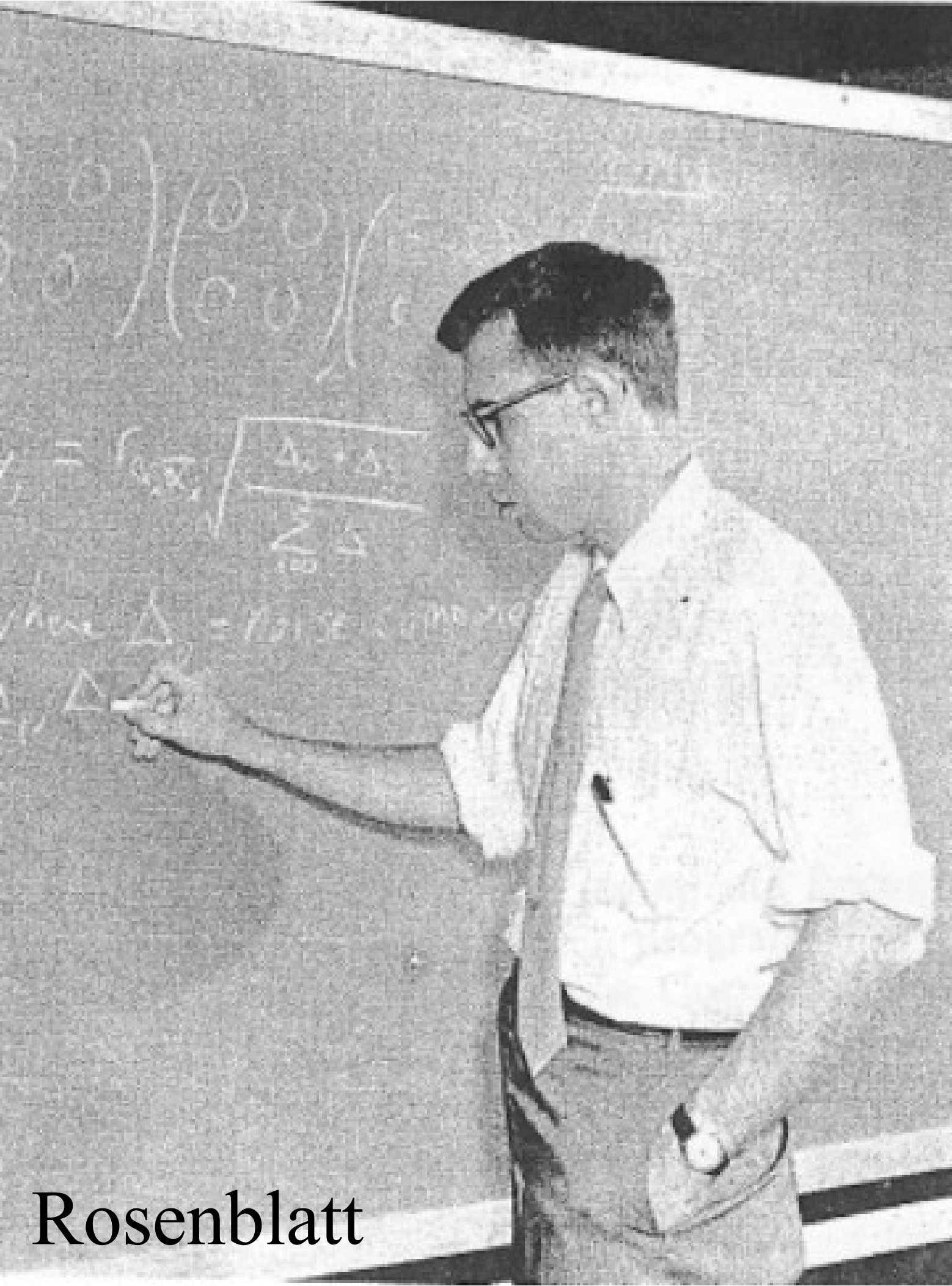
		x_2	
		0	1
x_1	0	0	1
	1	1	0

XOR

A brief history of Neural Networks

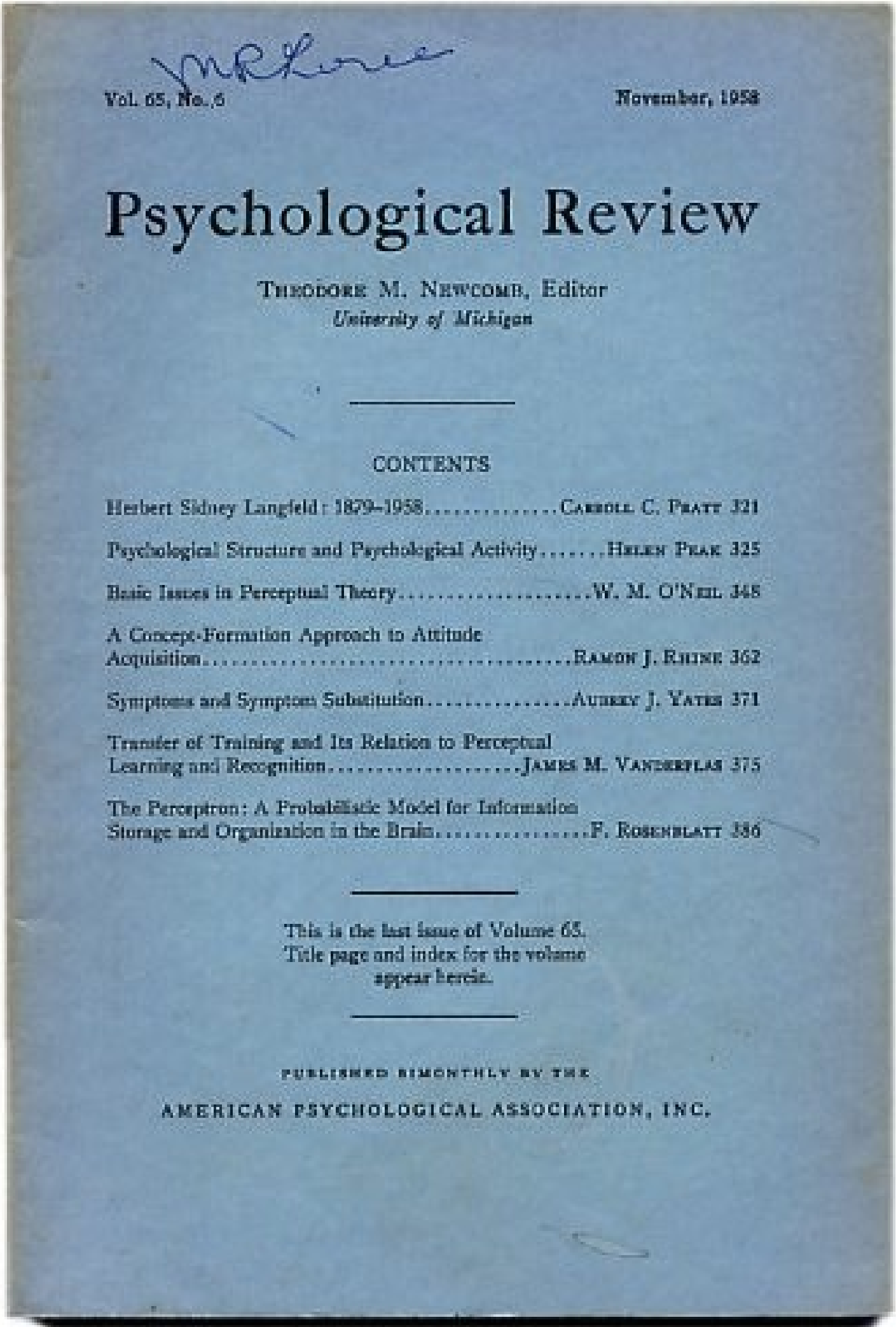


Perceptrons, 1958

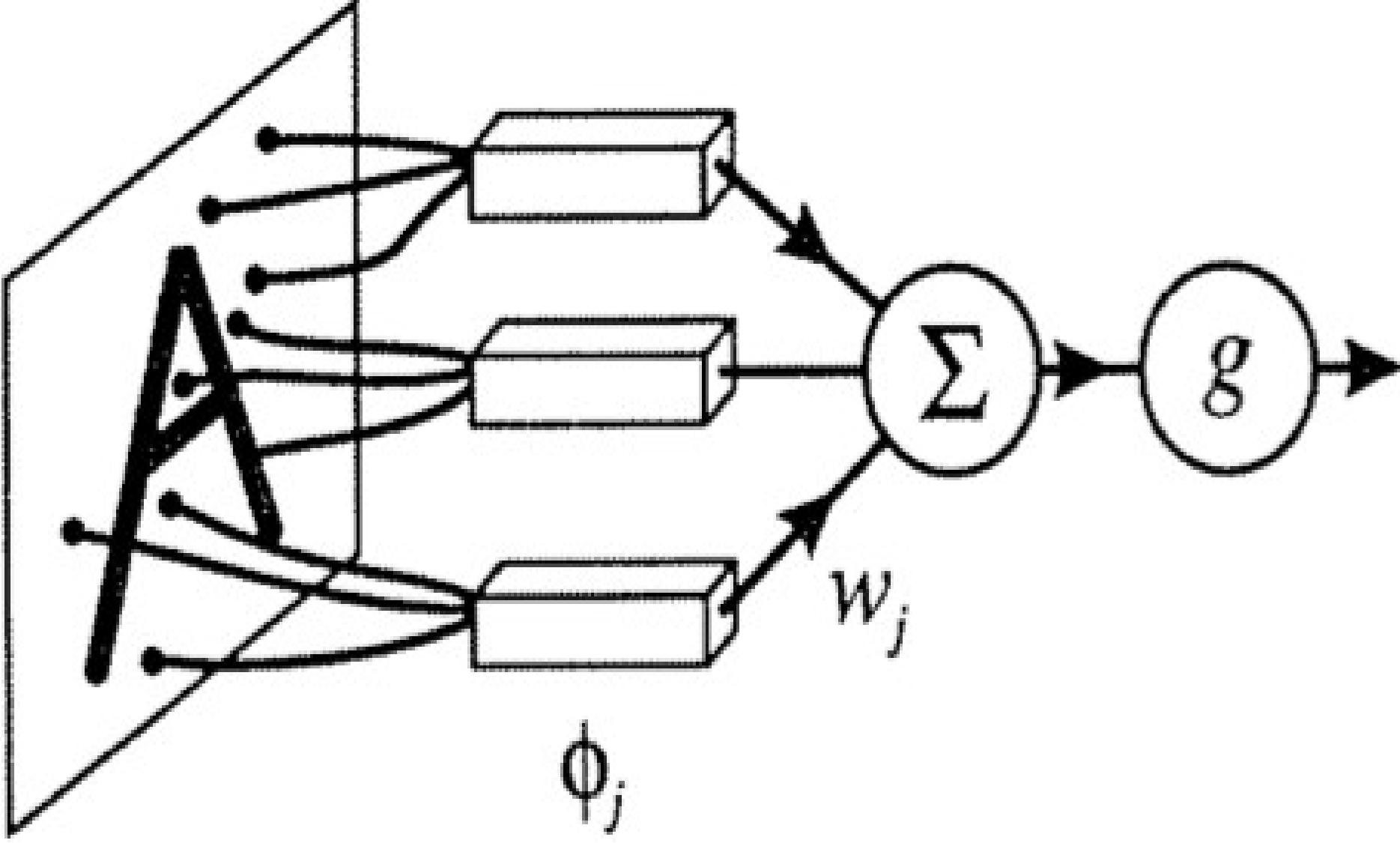


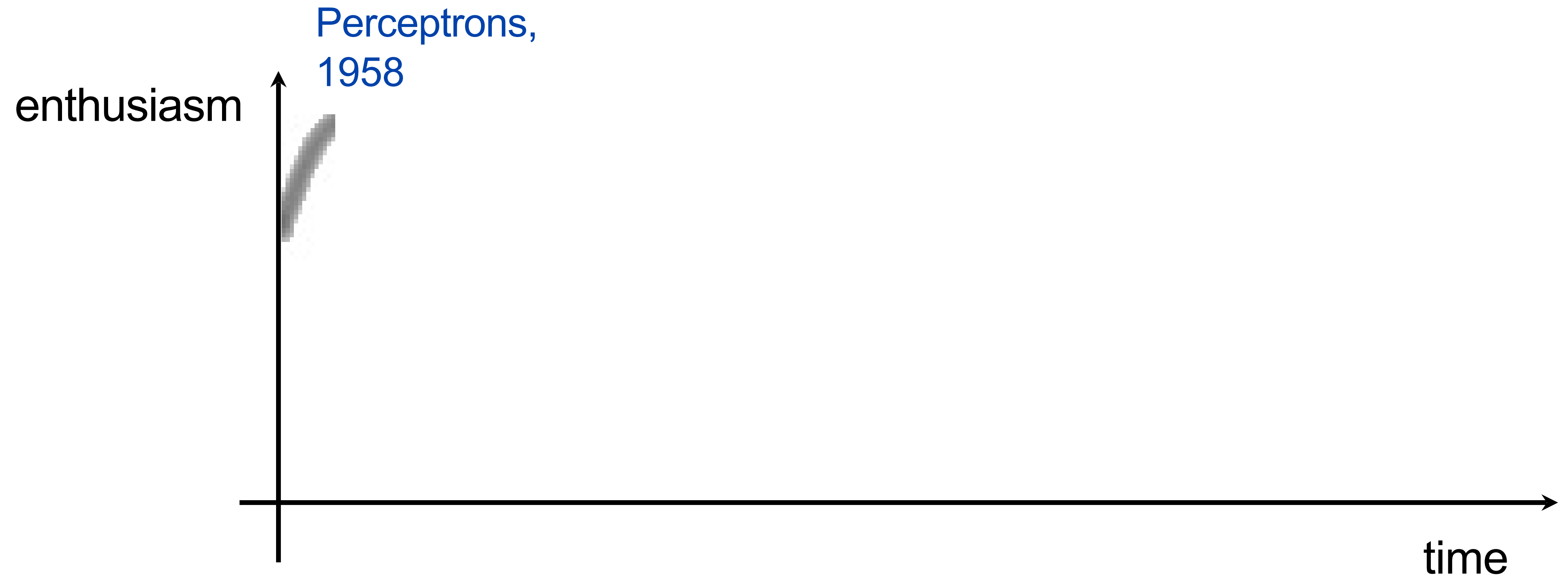
Rosenblatt

http://www.ecse.rpi.edu/homepages/nagy/PDF_chrono/2011_Nagy_Pace_FR.pdf. Photo by George Nagy

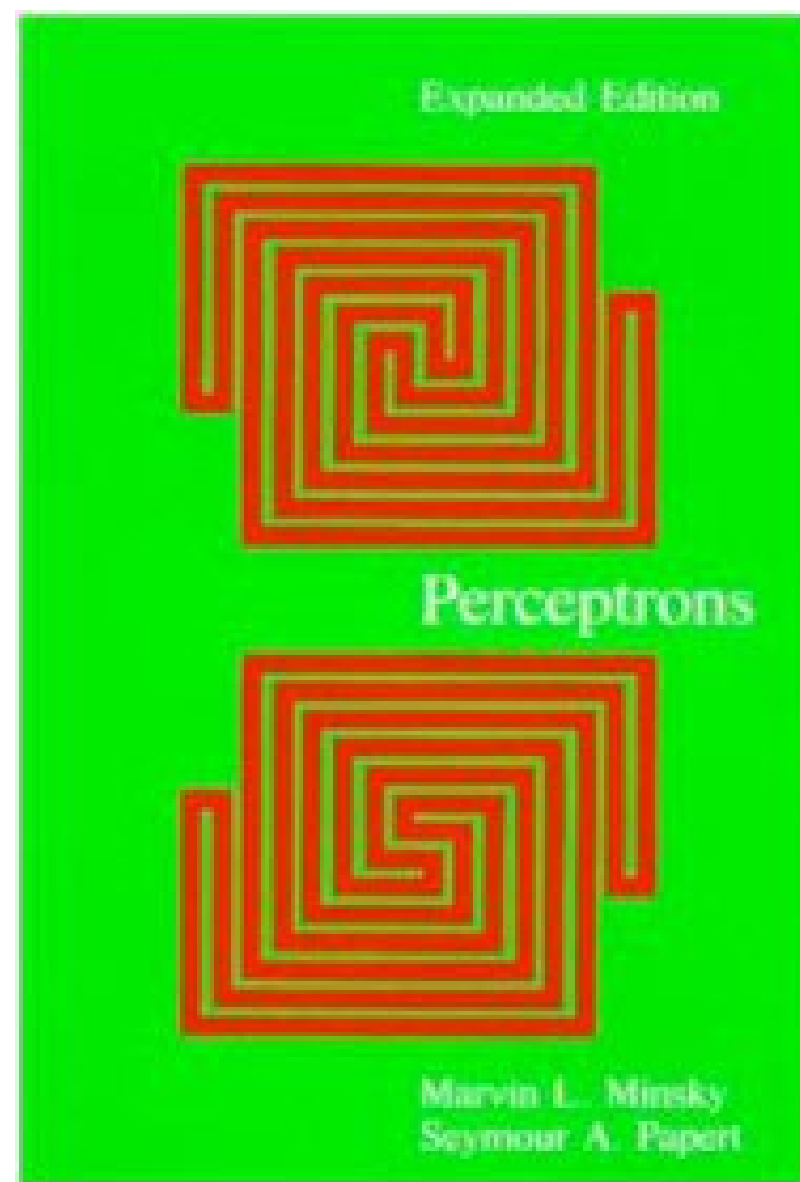


<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.335.3398&rep=rep1&type=pdf>





Minsky and Papert, Perceptrons, 1972



FOR BUYING OPTIONS, START HERE

Select Shipping Destination

Paperback | \$35.00 Short | £24.95 |
ISBN: 9780262631112 | 308 pp. | 6 x
8.9 in | December 1987

Perceptrons, expanded edition

An Introduction to Computational Geometry

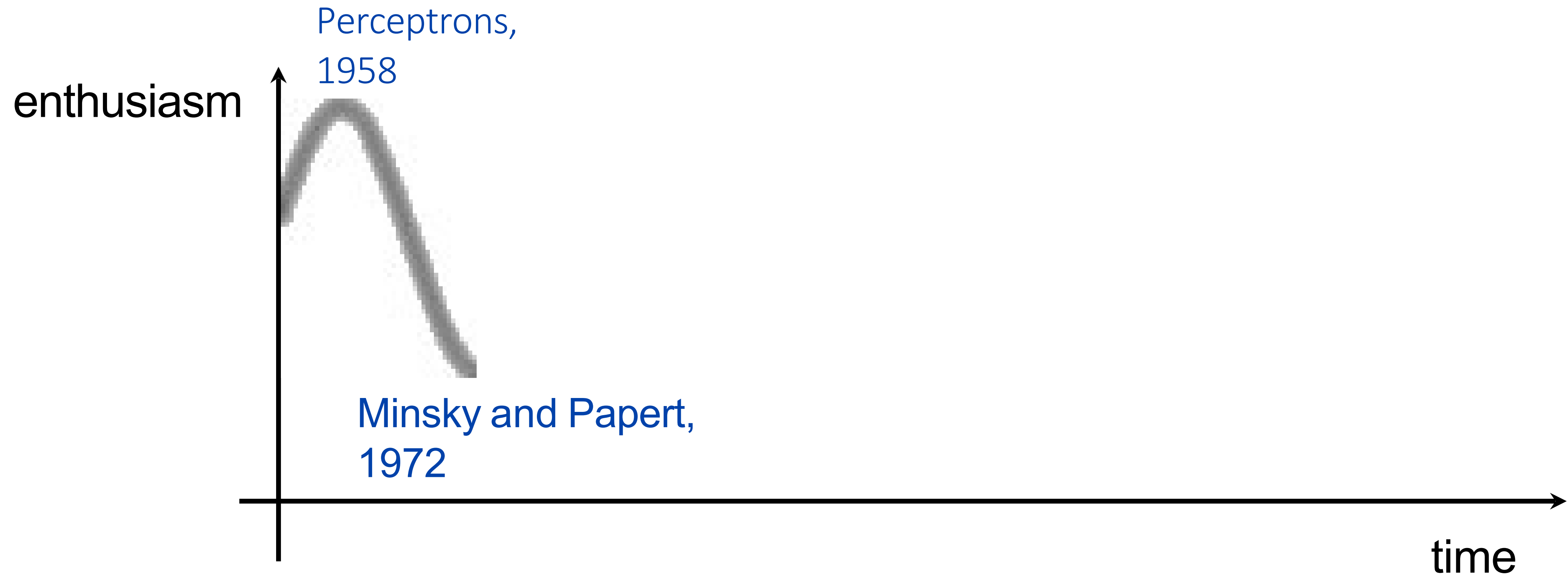
By [Marvin Minsky](#) and [Seymour A. Papert](#)

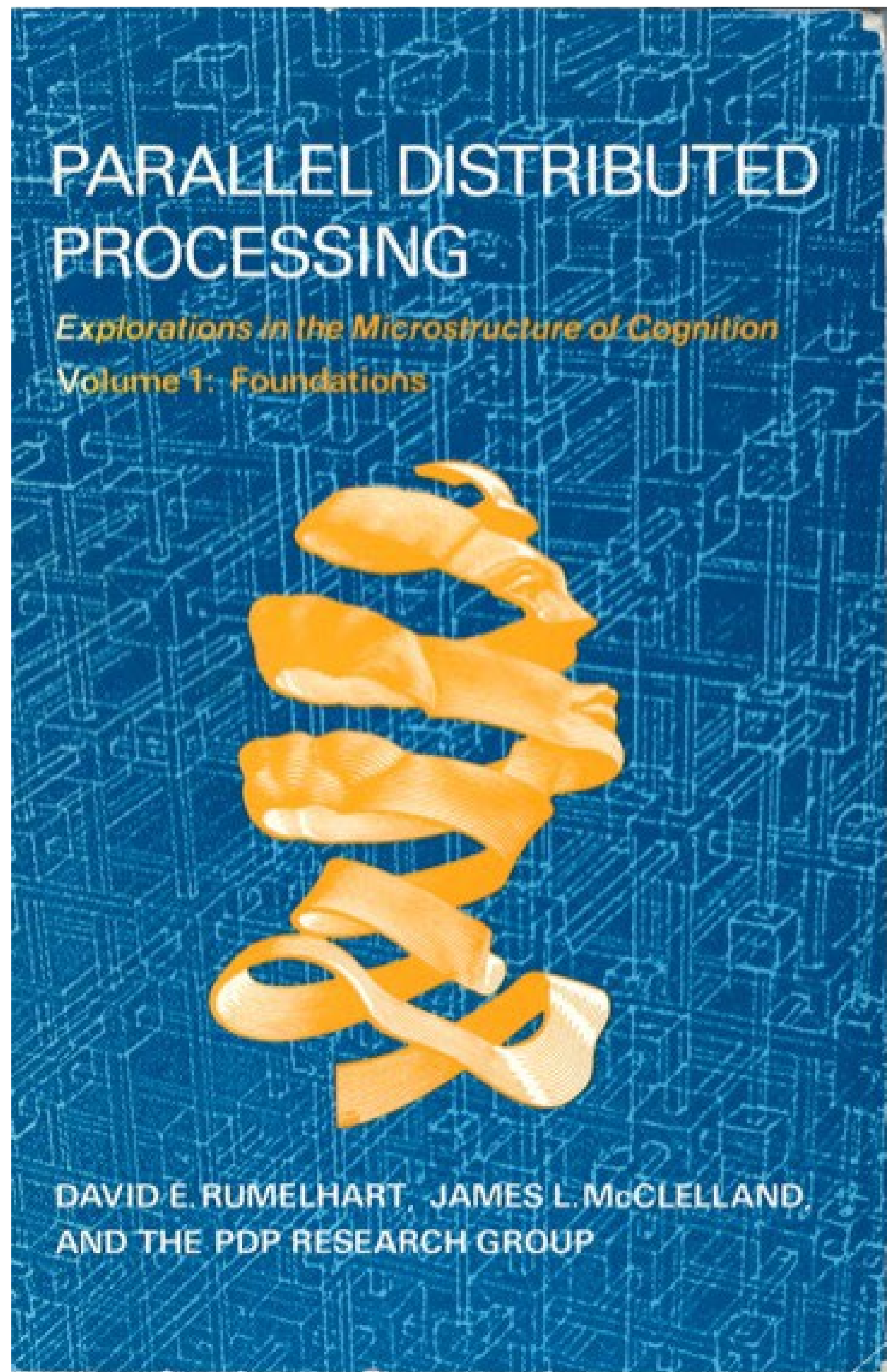
Overview

Perceptrons - the first systematic study of parallelism in computation - has remained a classical work on threshold automata networks for nearly two decades. It marked a historical turn in artificial intelligence, and it is required reading for anyone who wants to understand the connectionist counterrevolution that is going on today.

Artificial-intelligence research, which for a time concentrated on the programming of ton Neumann computers, is swinging back to the idea that intelligence might emerge from the activity of networks of neuronlike entities. Minsky and Papert's book was the first example of a mathematical analysis carried far enough to show the exact limitations of a class of computing machines that could seriously be considered as models of the brain. Now the new developments in mathematical tools, the recent interest of physicists in the theory of disordered matter, the new insights into and psychological models of how the brain works, and the evolution of fast computers that can simulate networks of automata have given *Perceptrons* new importance.

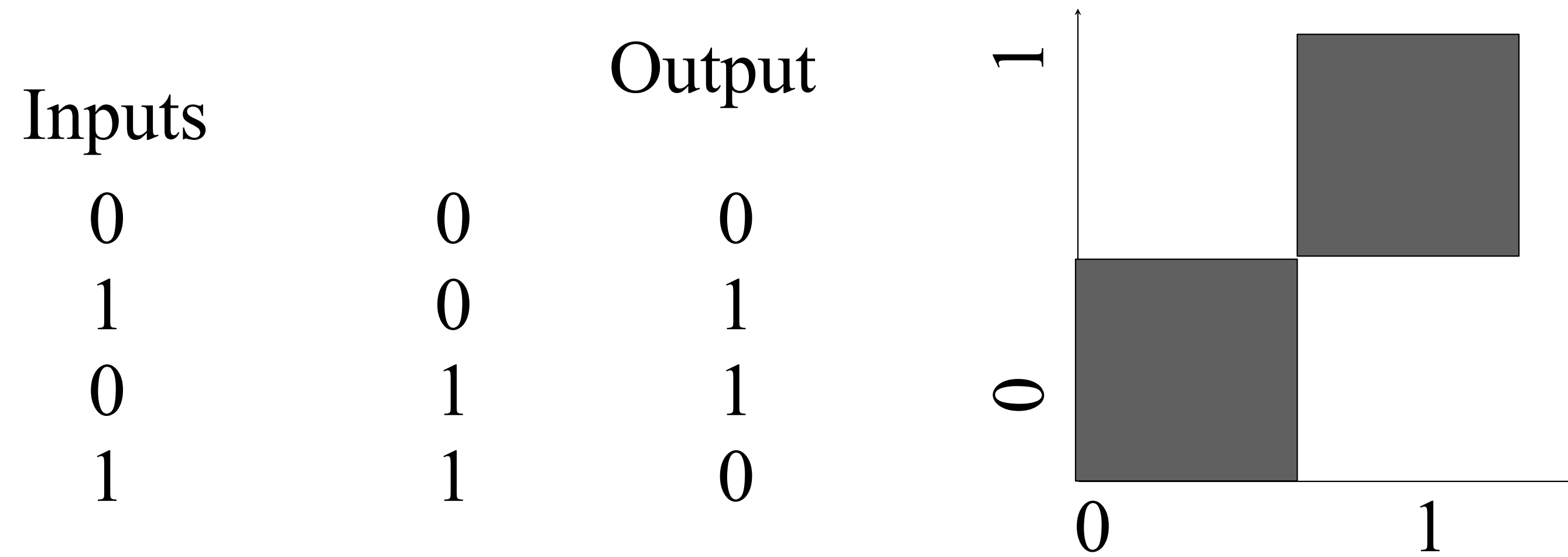
Witnessing the swing of the intellectual pendulum, Minsky and Papert have added a new chapter in which they discuss the current state of parallel computers, review developments since the appearance of the 1972 edition, and identify new research directions related to connectionism. They note a central theoretical challenge facing connectionism: the challenge to reach a deeper understanding of how "objects" or "agents" with individuality can emerge in a network. Progress in this area would link connectionism with what the authors have called "society theories of mind."



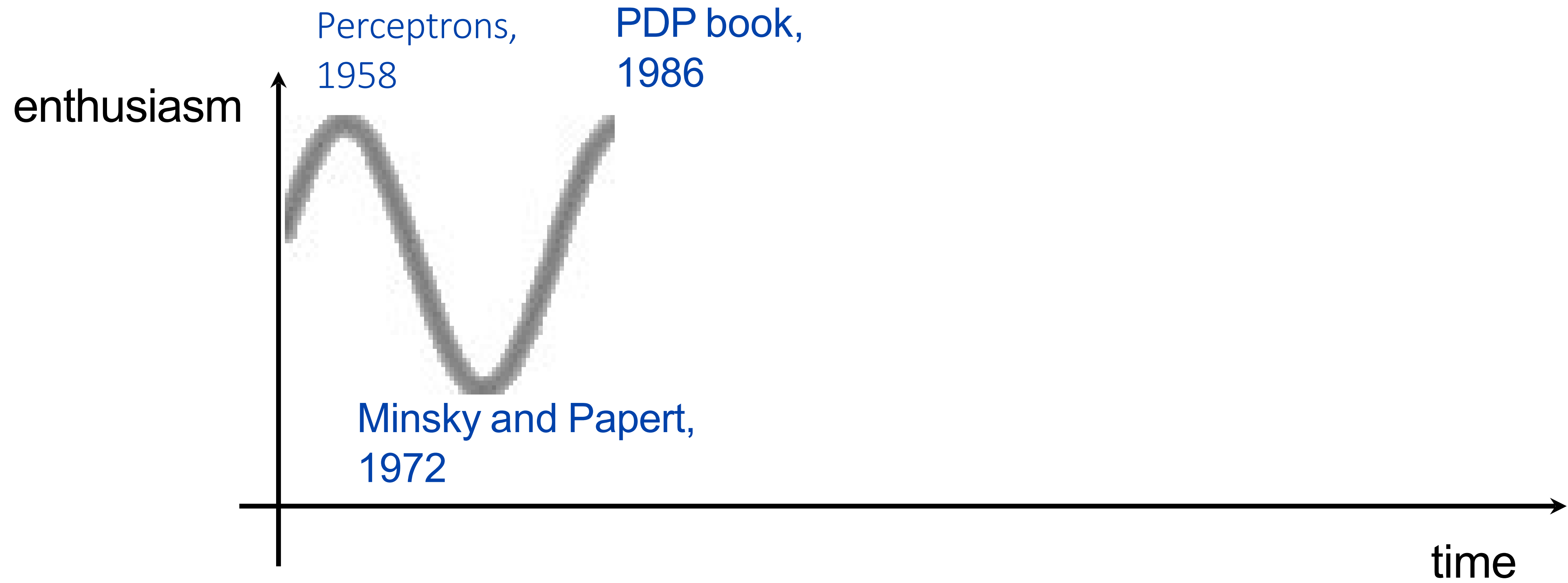


Parallel Distributed Processing (PDP), 1986

XOR problem



PDP authors pointed to the backpropagation algorithm as a breakthrough, allowing multi-layer neural networks to be trained. Among the functions that a multi-layer network can represent but a single-layer network cannot: the XOR function.



LeCun convolutional neural networks



PROC. OF THE IEEE, NOVEMBER 1998

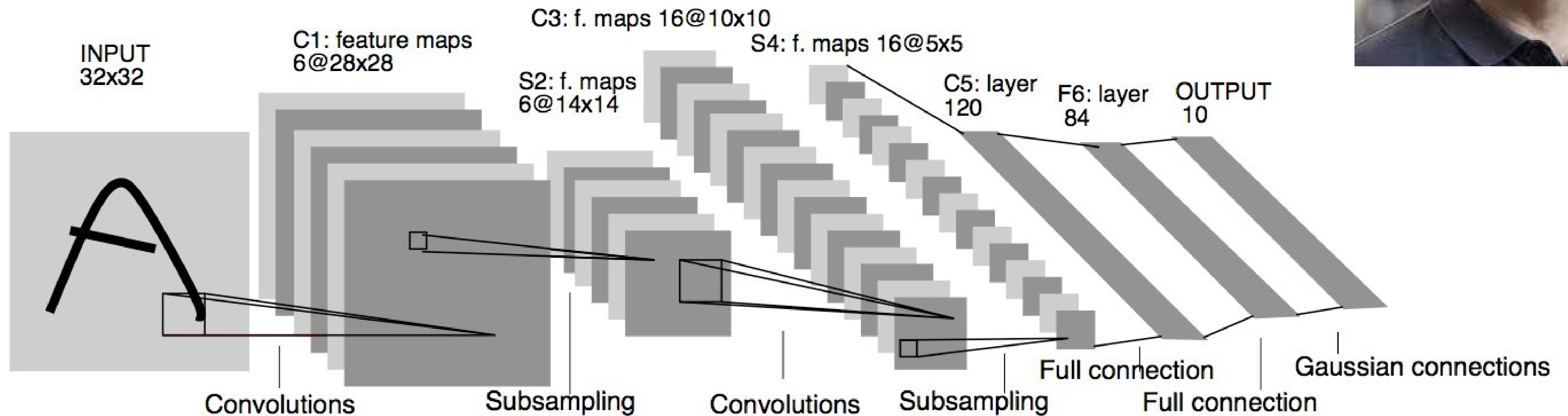


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Demos:

<http://yann.lecun.com/exdb/lenet/index.html>



Yann LeCun

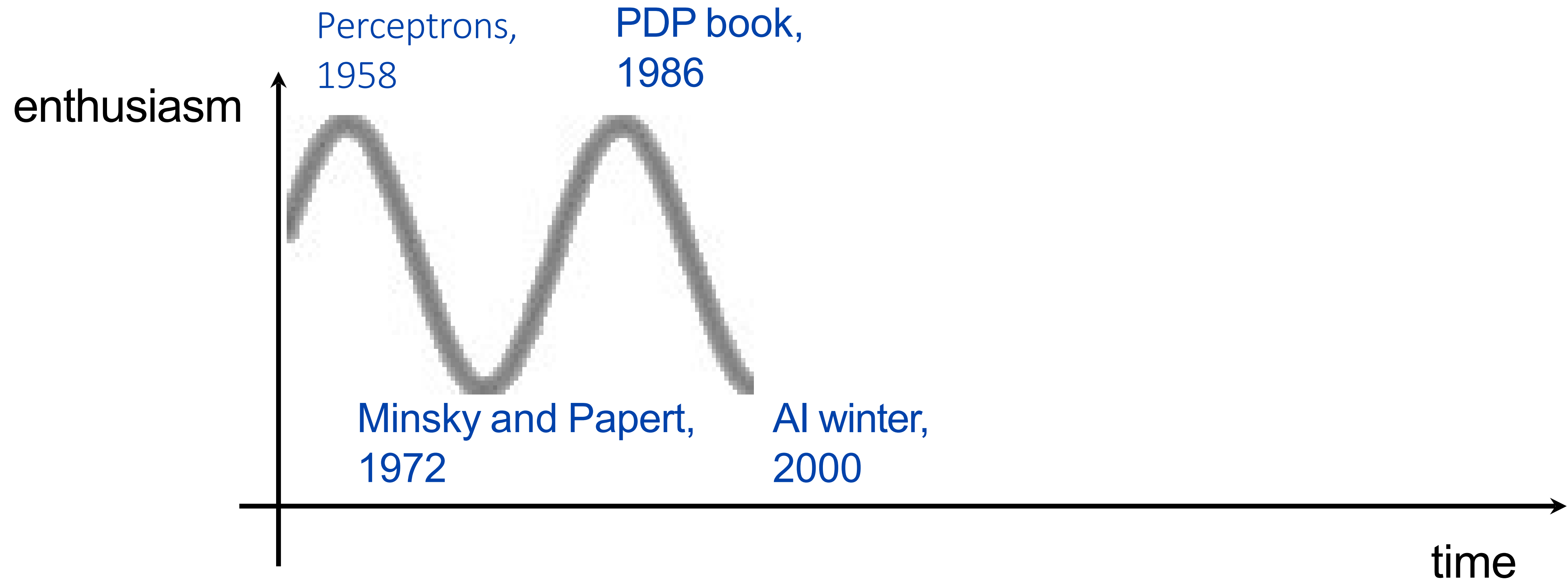
Was at Bell Labs when
this video was recorded

Now
Prof @ NYU
Chief Scientist @ Meta

Turing Award 2018
(shared with Hinton and
Bengio)

Neural Information Processing Systems 2000

- Neural Information Processing Systems is the premier conference on machine learning.
- Evolved from an interdisciplinary conference to a machine learning conference.
- For the 2000 conference:
 - title words predictive of paper acceptance: “Belief Propagation” and “Gaussian”.
 - title words predictive of paper rejection: “Neural” and “Network”.



ImageNet: First (?) large-scale computer vision dataset

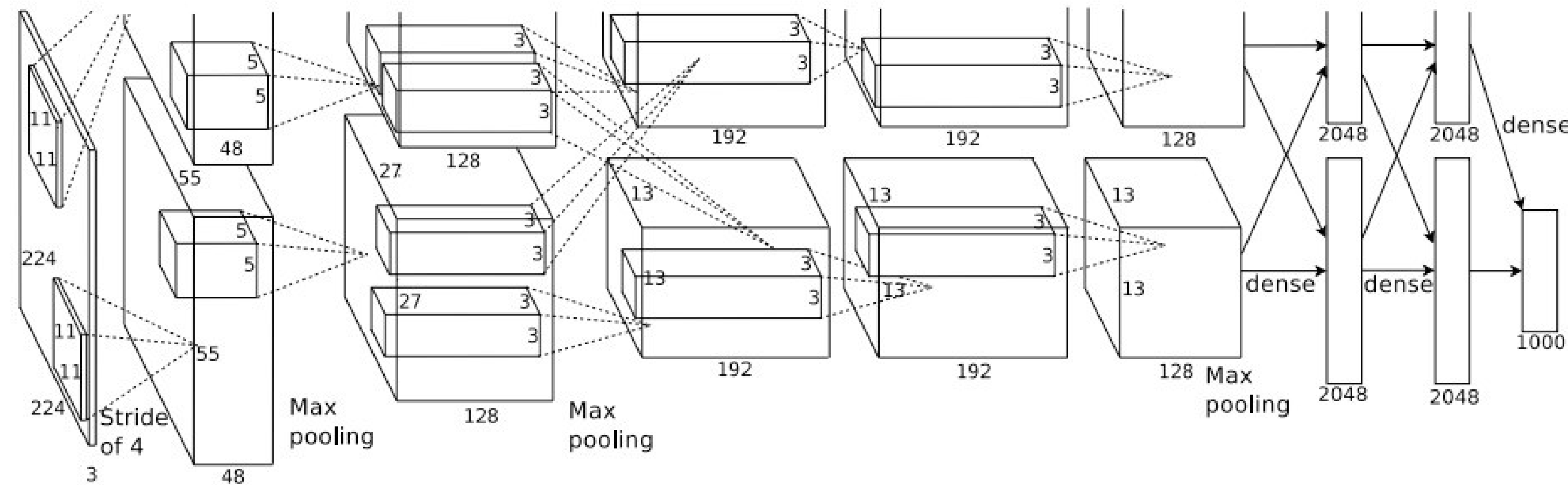


- Millions of images; 1000 categories
- **PI: Fei-Fei Li**
 - Then: Prof, Princeton
 - Now: Prof, Stanford
- 2019 Longuet-Higgins Prize
 - Some argued that Li deserved the 2018 Turing Award along with Hinton, LeCun, Bengio
 - Their work could not have been empirically tested without ImageNet!



Krizhevsky, Sutskever, and Hinton, NeurIPS 2012

“AlexNet”

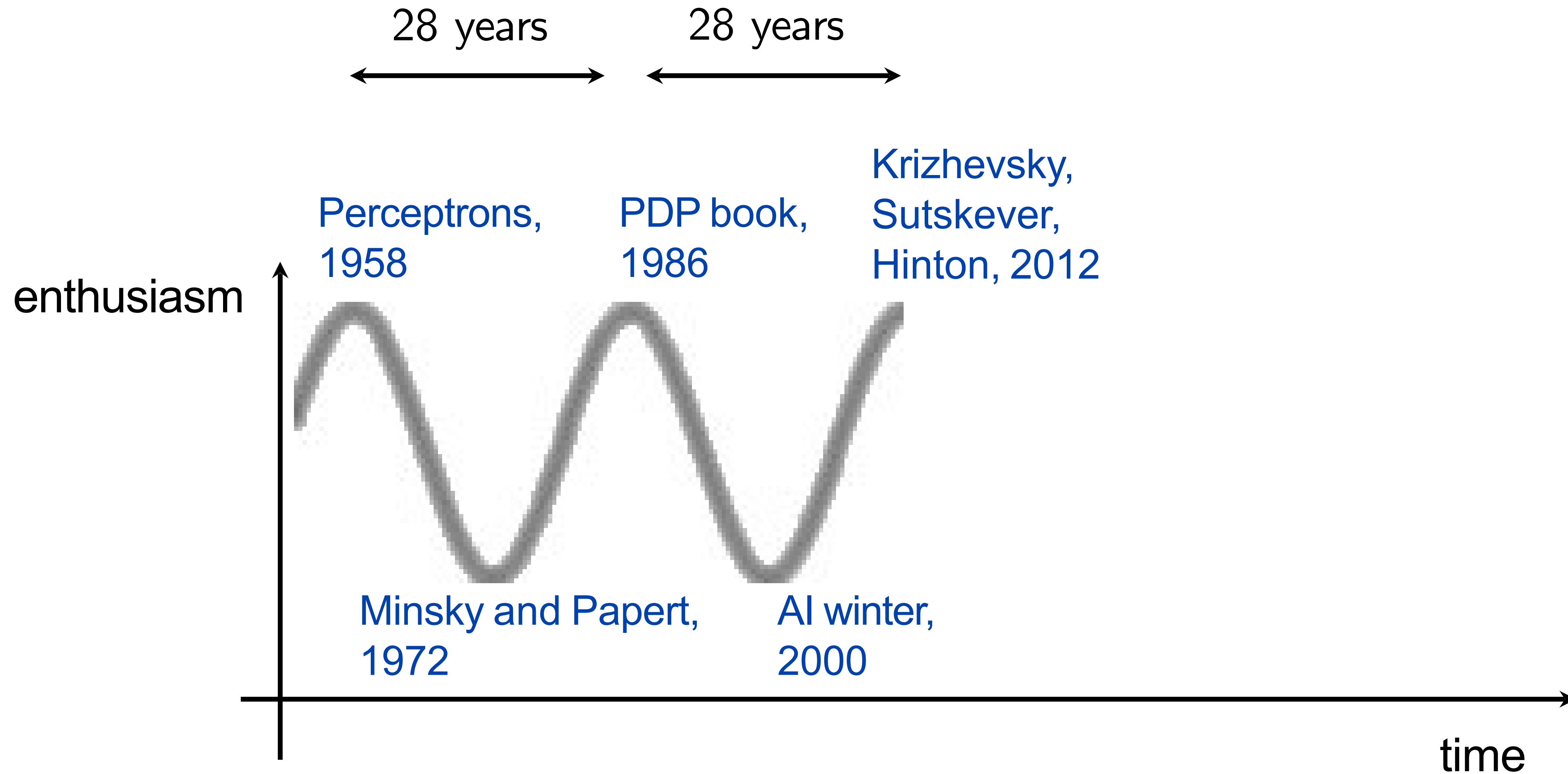


Got all the “pieces” right, e.g.,

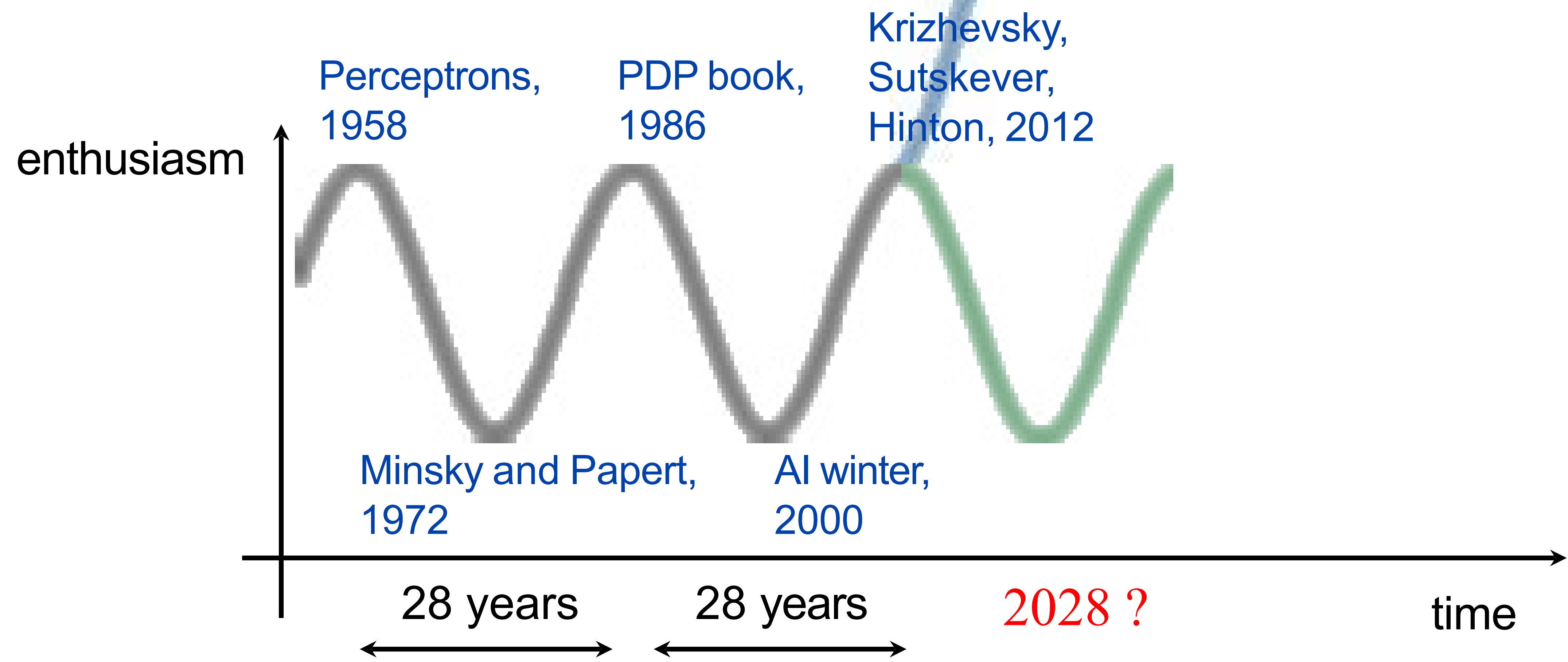
- **Trained on ImageNet**
- 8 layer architecture (for reference: today we have architectures with 100+ layers)
- Allowed for multi-GP training

Krizhevsky, Sutskever, and Hinton, NeurIPS 2012

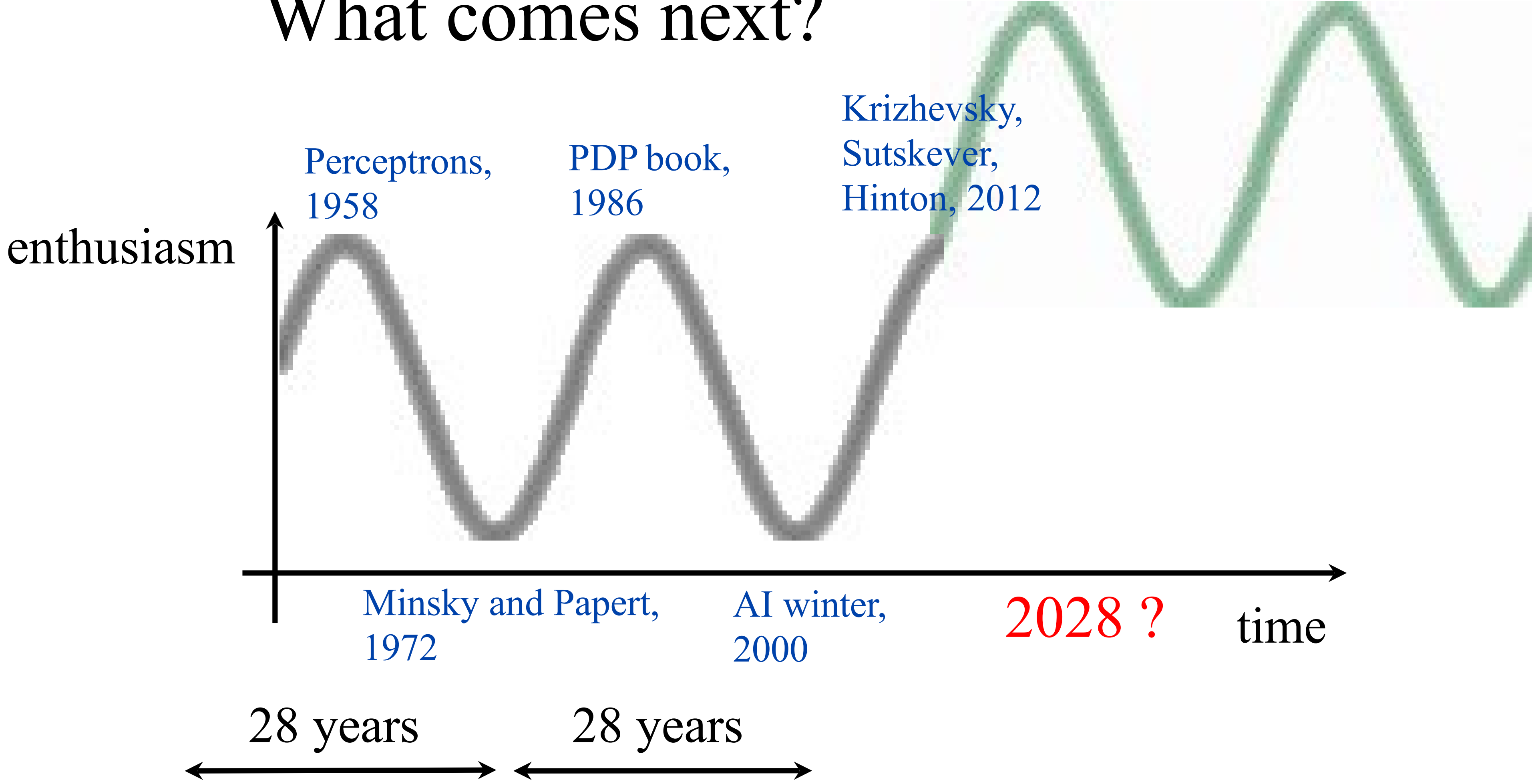




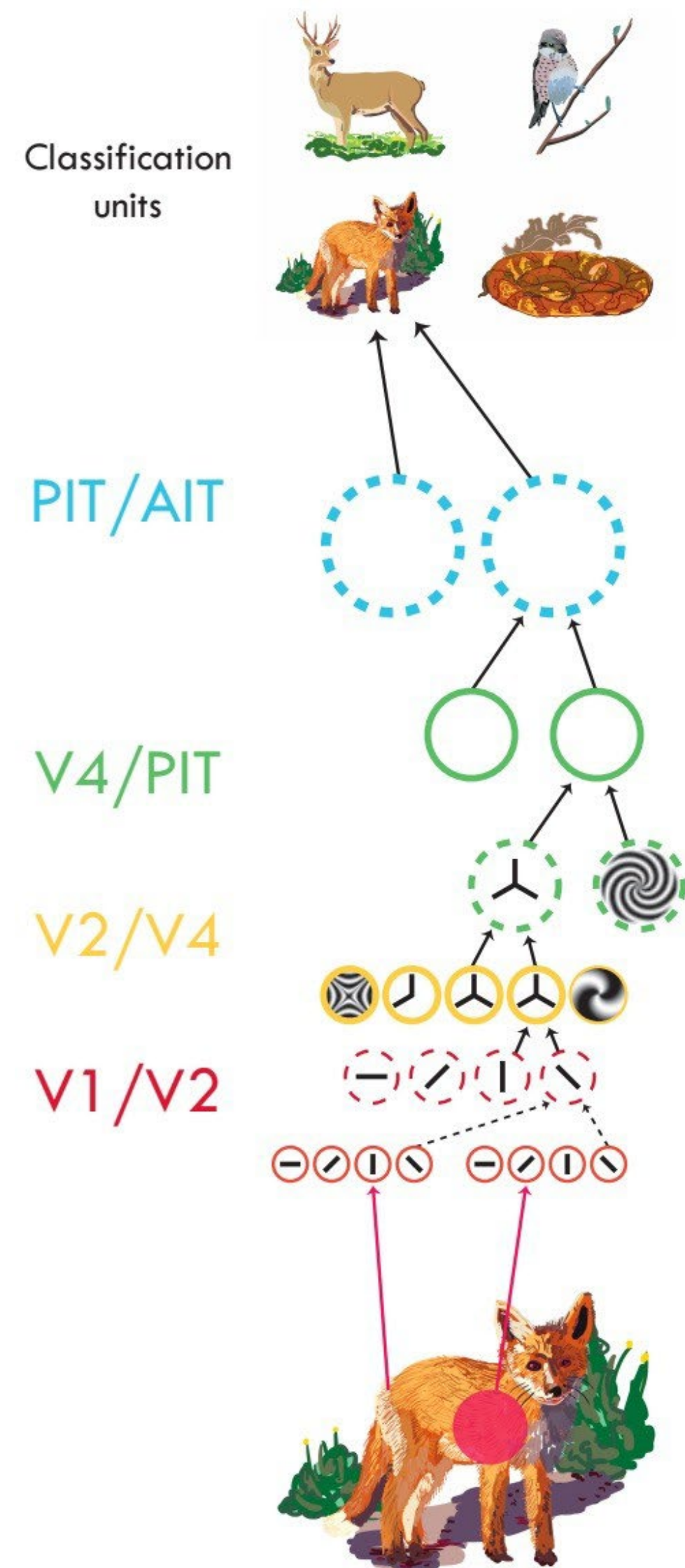
What comes next?



What comes next?



Inspiration: Hierarchical Representations



Classification units

PIT/AIT

V4/PIT

V2/V4

V1/V2

[Serre, 2014]

Best to treat as *inspiration*.

The neural nets we'll talk about aren't very biologically plausible.

Object recognition

Pixel 1



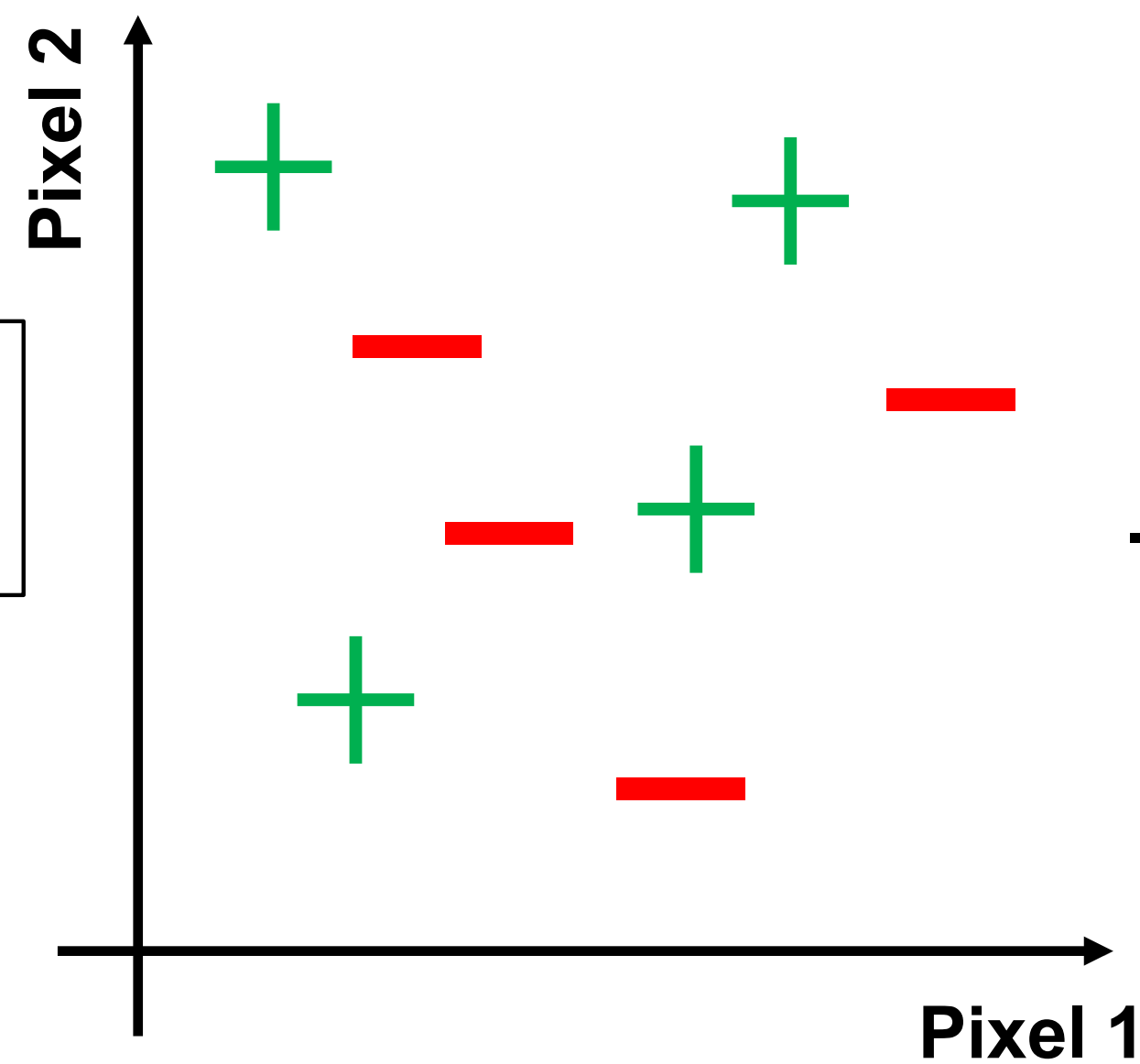
Neural Network

Is dog?

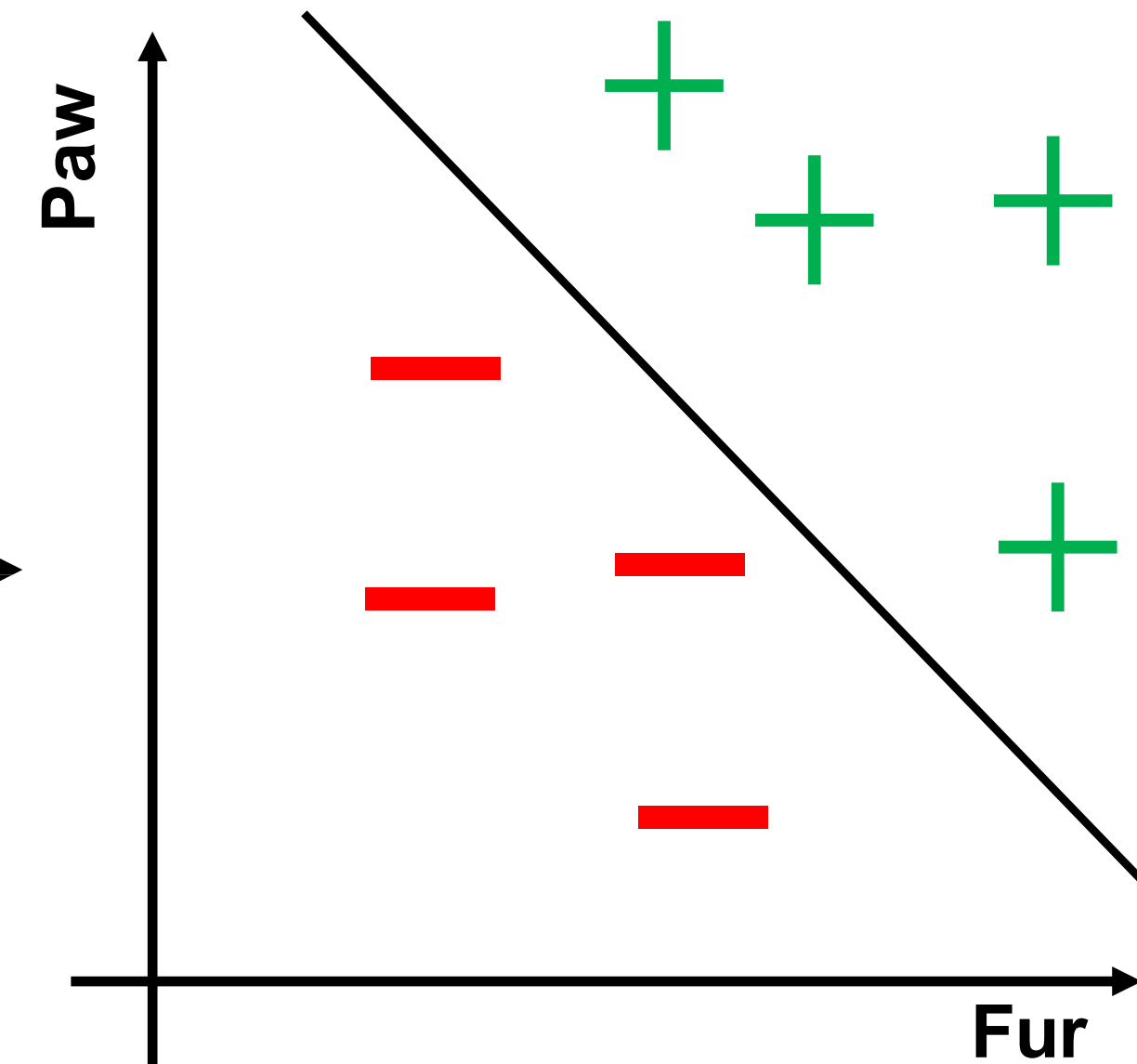
Pixel 2

Input Space

Feature Space



$f(x)$



+ Dog
- Not dog

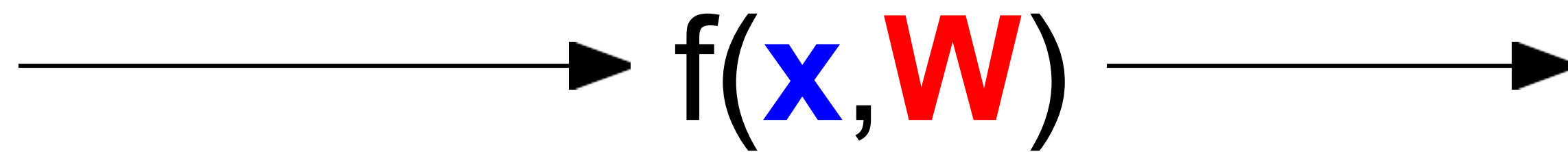
Goal: automatically learn a function that maps data from the input space to a feature space, i.e., "feature learning", rather than use hand-crafted features

Parametric Approach

Image



Array of **32x32x3** numbers
(3072 numbers total)



10 numbers giving
class scores

↑
W

parameters
or weights

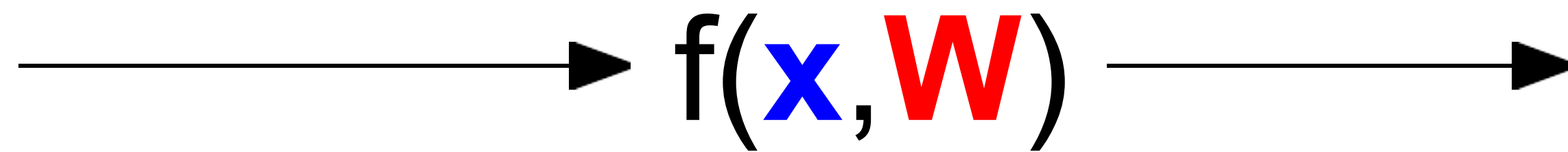
Parametric Approach: Linear Classifier

$$f(x, W) = Wx$$

Image



Array of **32x32x3** numbers
(3072 numbers total)



10 numbers giving
class scores



W

parameters
or weights

Parametric Approach: Linear Classifier

Image

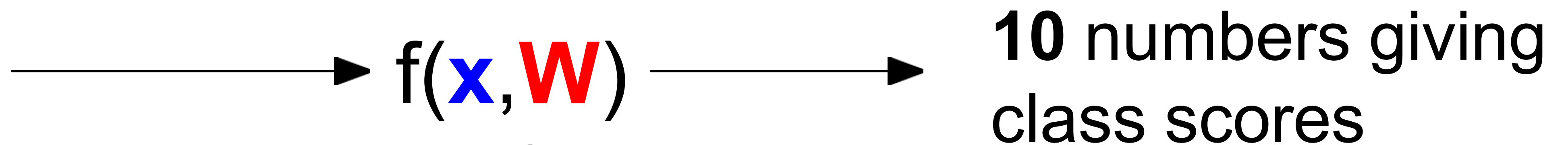


Array of $32 \times 32 \times 3$ numbers
(3072 numbers total)

$$f(x, W) = Wx$$

3072×1

10×1 10×3072



W

parameters
or weights

Parametric Approach: Linear Classifier

Image

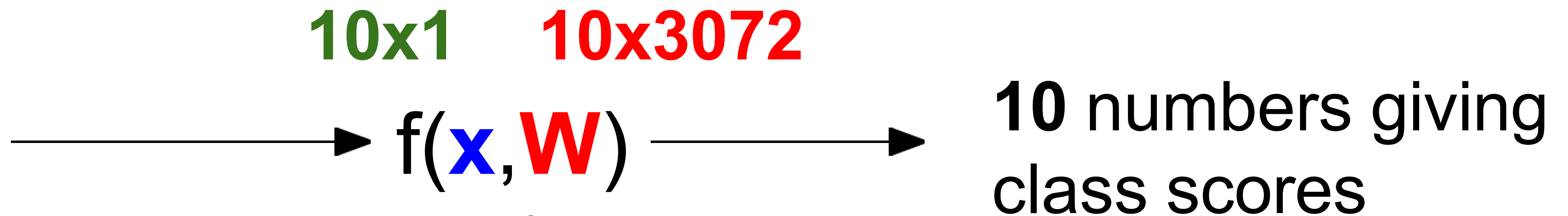


Array of $32 \times 32 \times 3$ numbers
(3072 numbers total)

$$f(x, W) = Wx + b$$

3072×1

10×1 10×3072 10×1



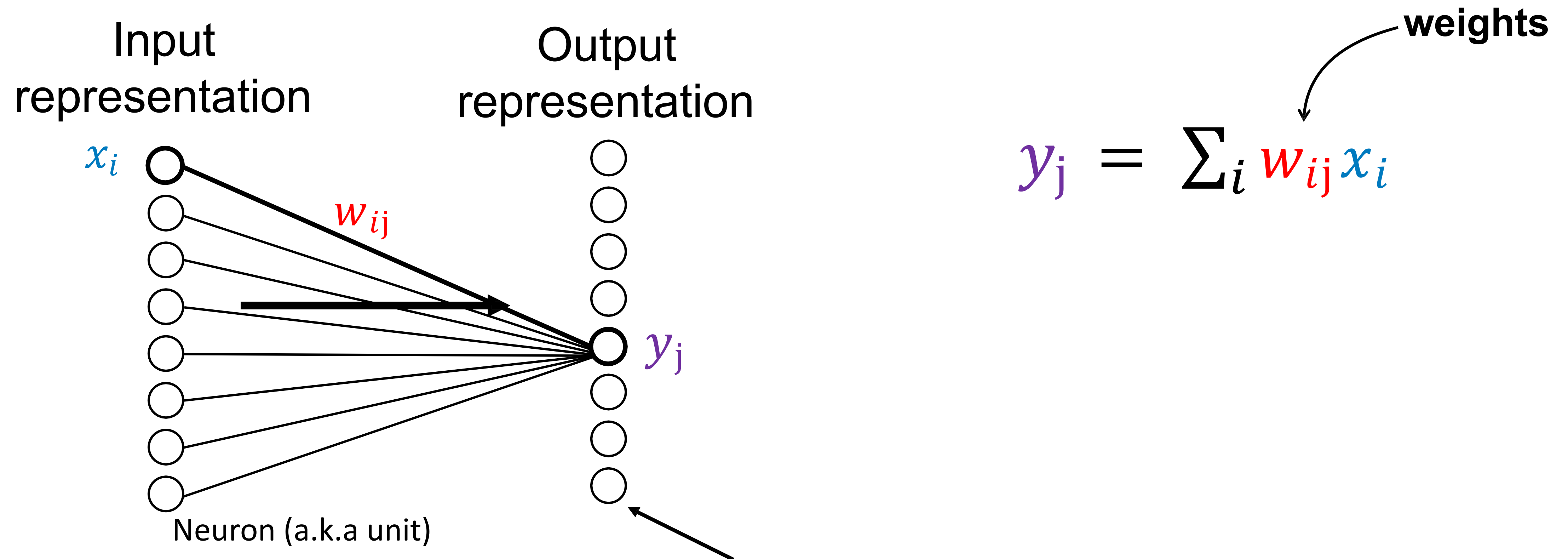
W

parameters
or weights

Computation in a neural net

Let's say we have some 1D input that we want to convert to some new feature space:

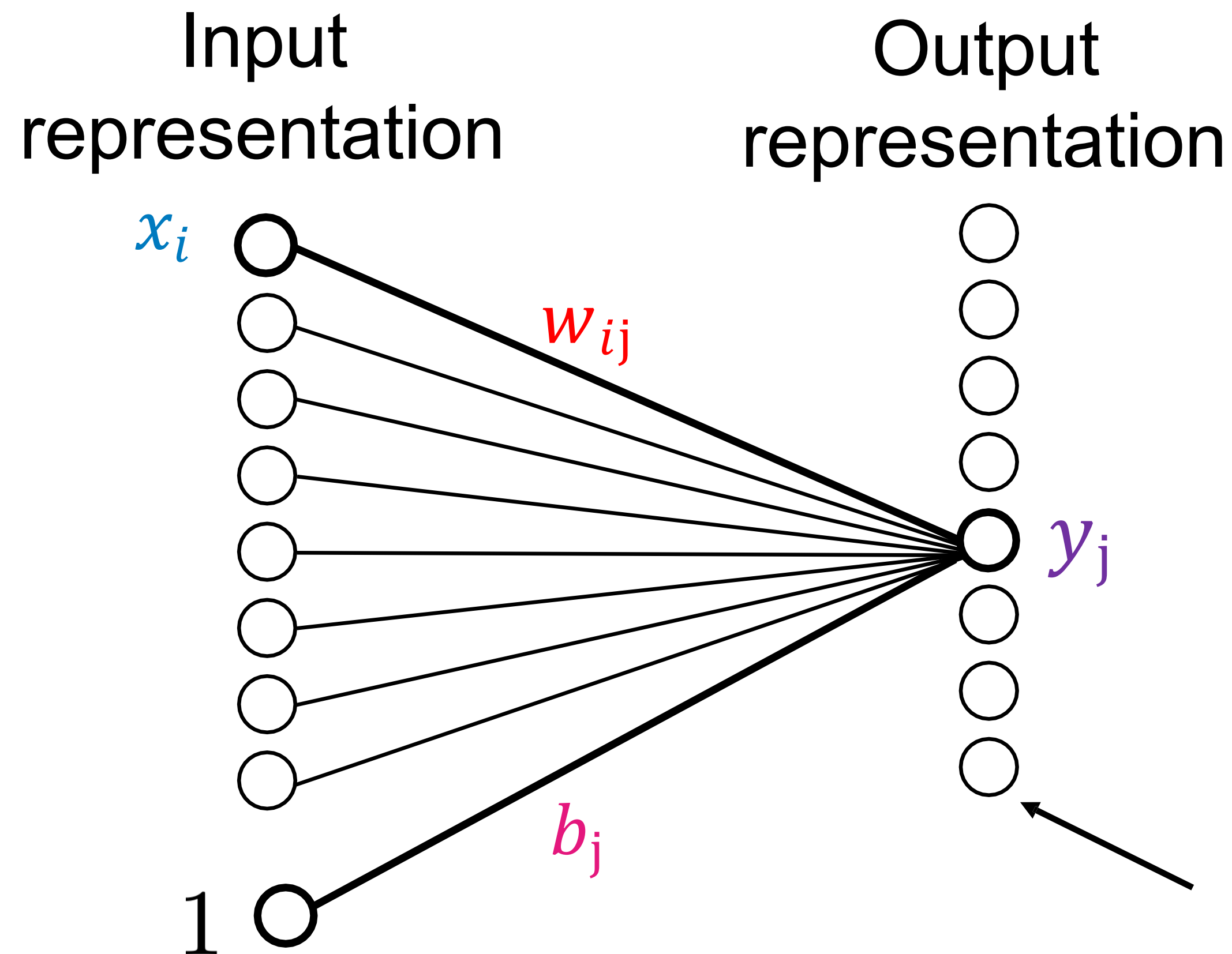
Linear layer



Computation in a neural net

Let's say we have some 1D input that we want to convert to some new feature space

Linear layer



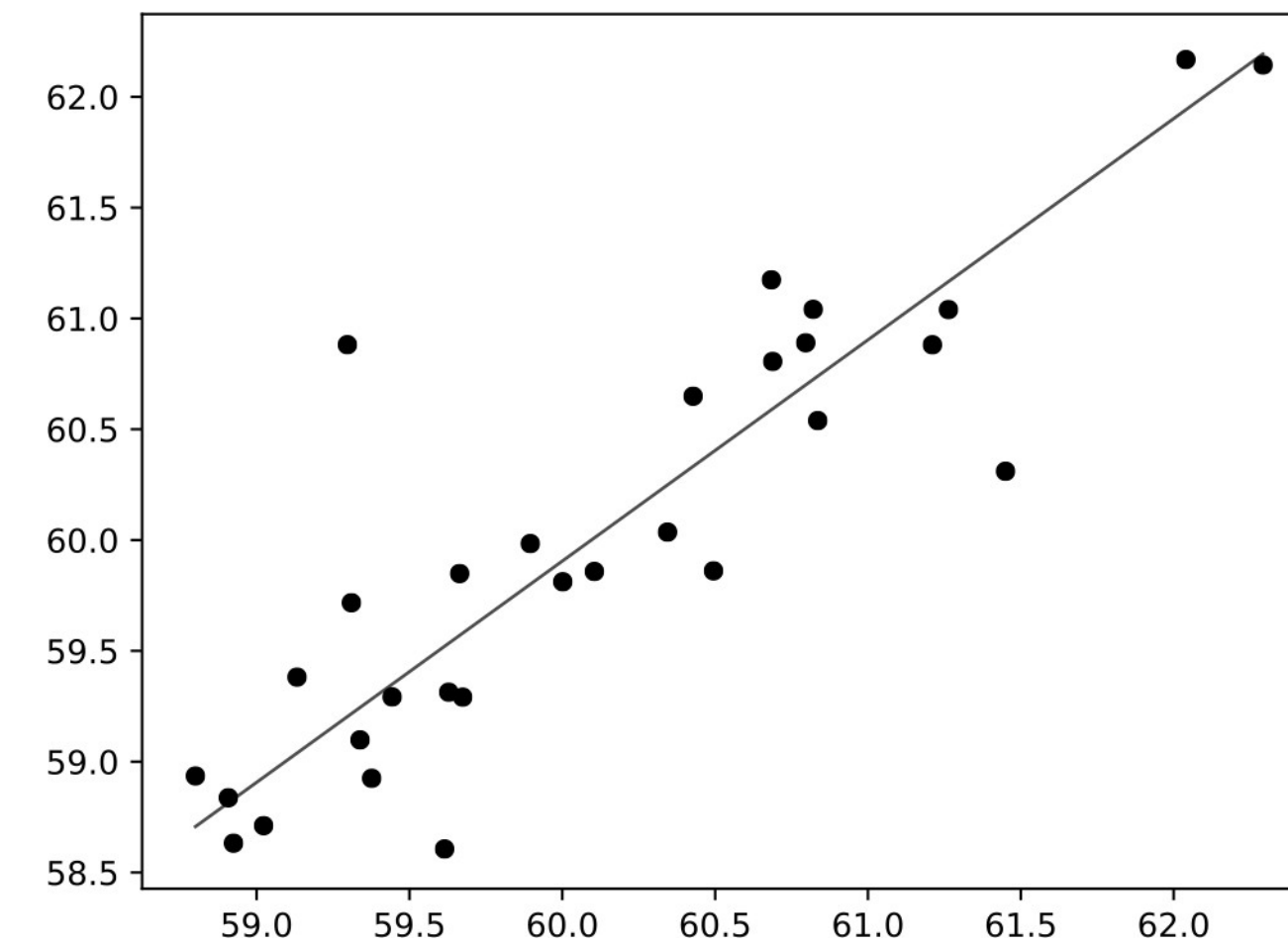
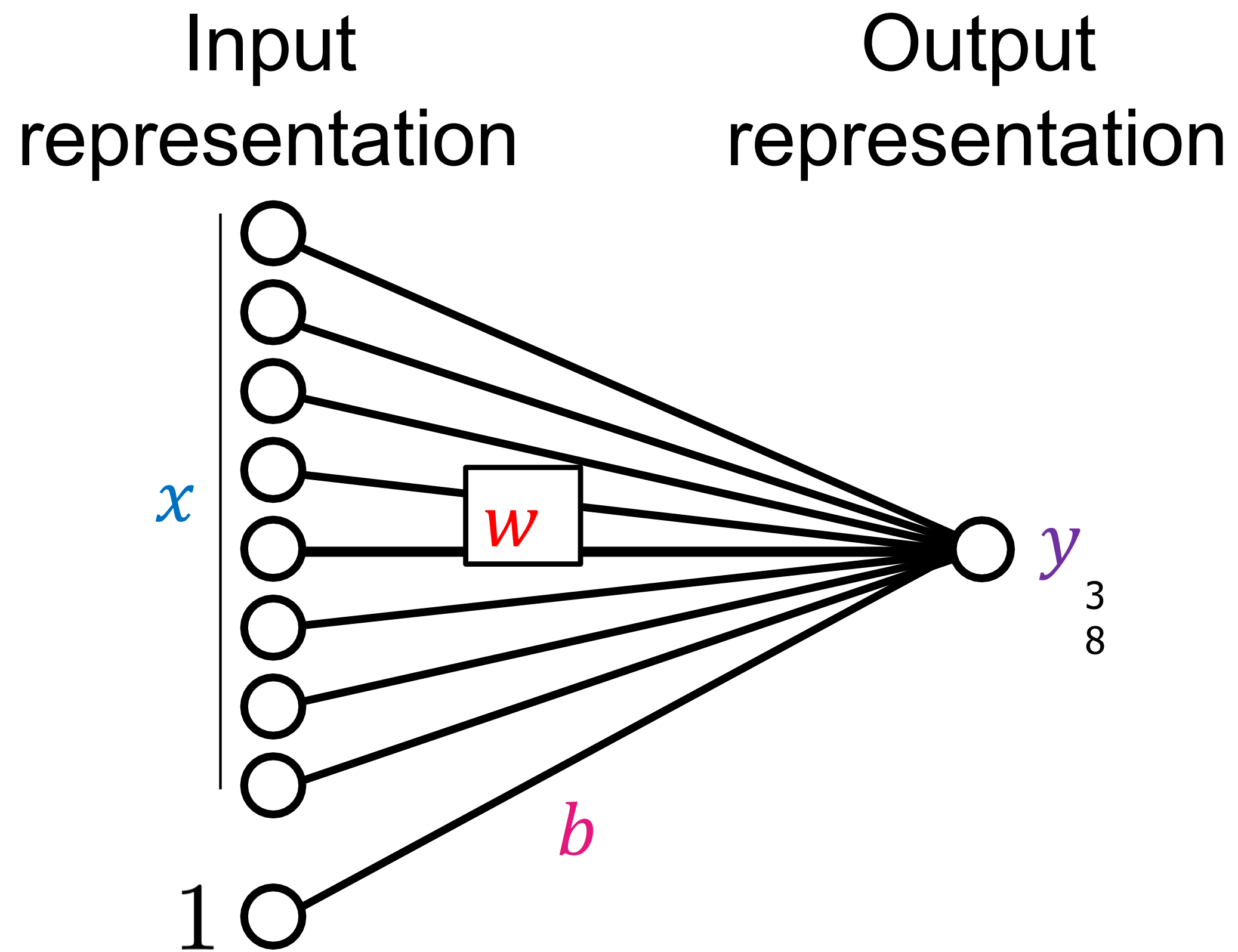
$$y_j = \sum_i w_{ij} x_i + b_j$$

weights

bias

Example: Linear Regression

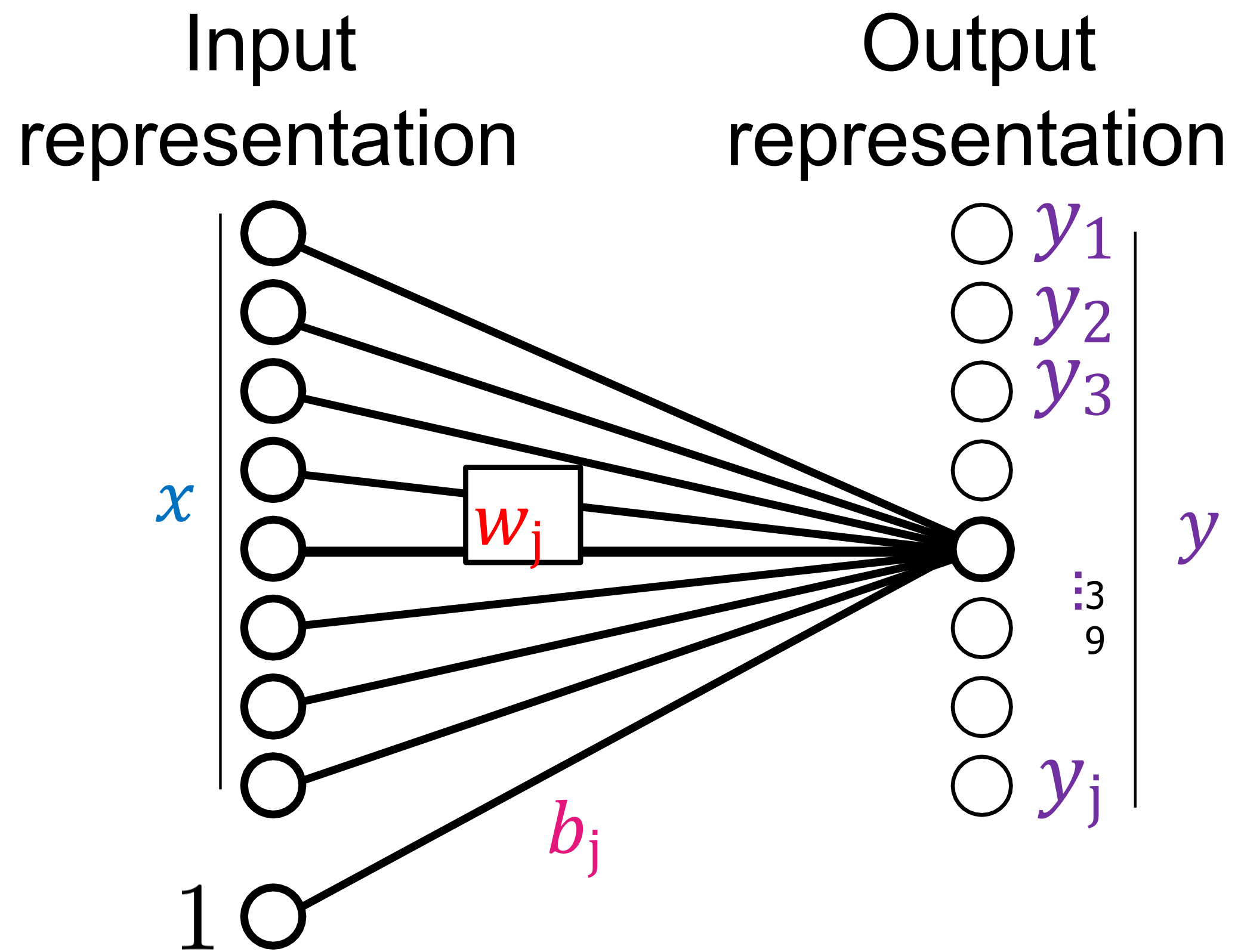
Linear layer



$$f_{\mathbf{w},b}(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b$$

Computation in a neural net – Full Layer

Linear layer



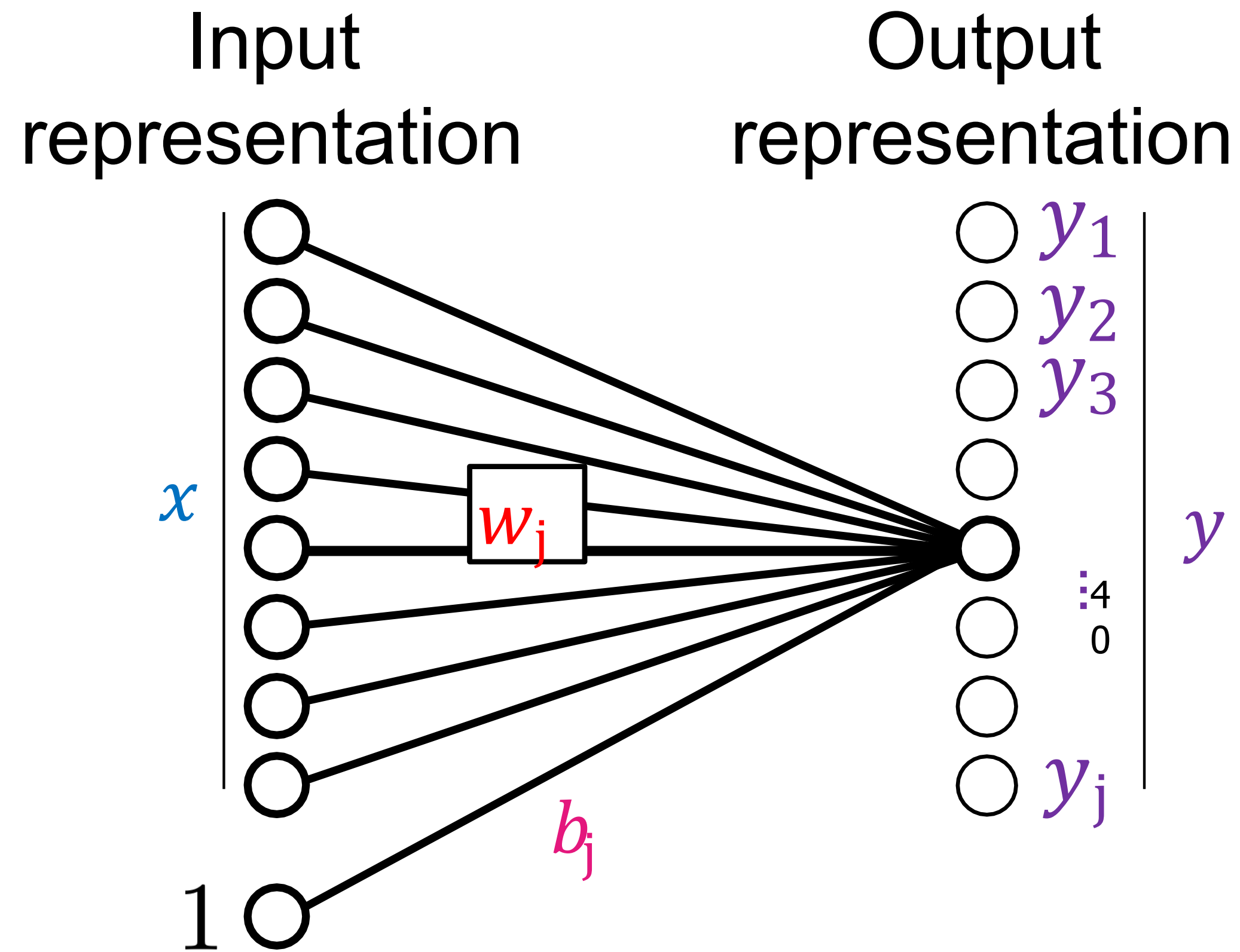
$$y = Wx + b$$

$$\begin{bmatrix} W_{11} & \cdots & W_{1n} \\ \vdots & \ddots & \vdots \\ W_{j1} & \cdots & W_{jn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_n \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \cdots \\ b_j \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_j \end{bmatrix}$$

parameters of the model: $\theta = \{W, b\}$

Computation in a neural net – Full Layer

Linear layer



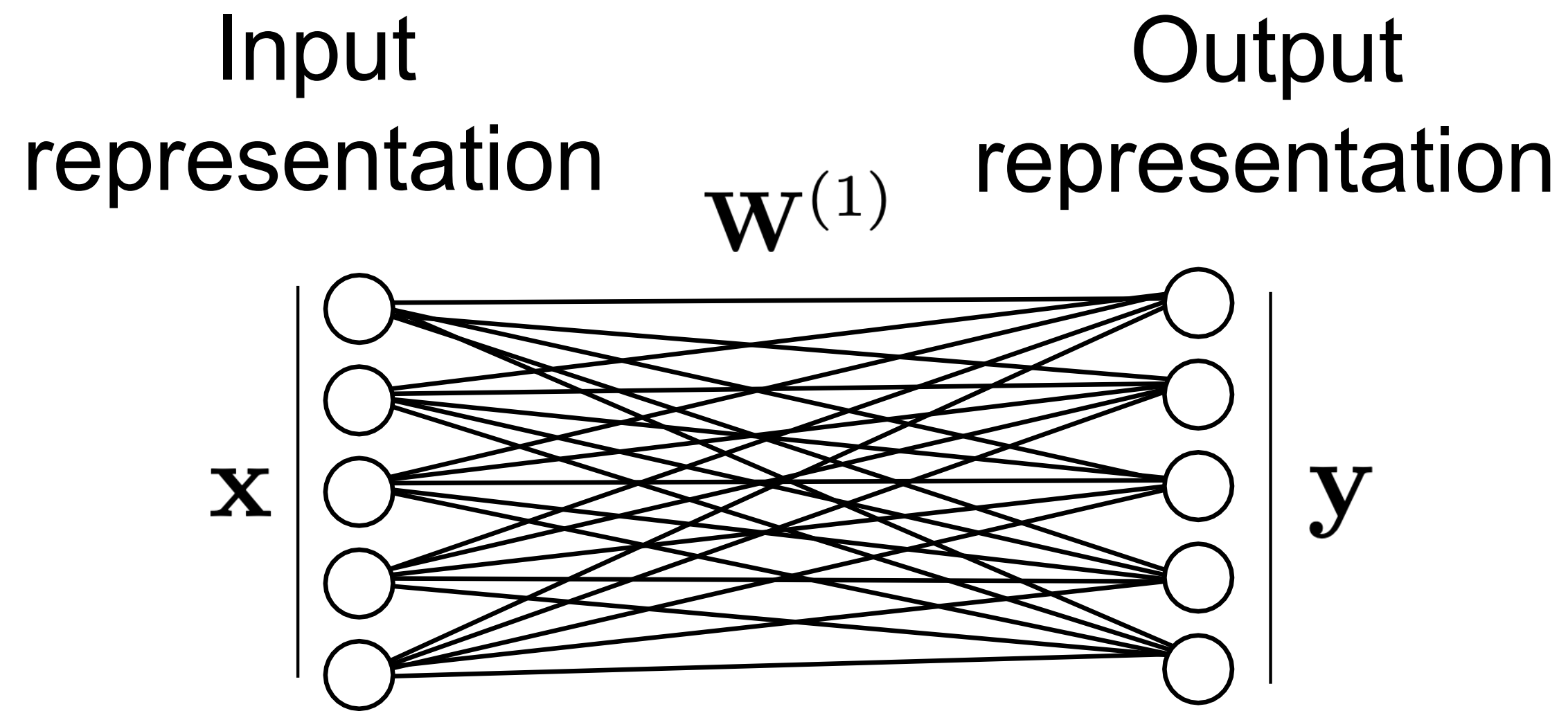
Full layer

$$y = Wx + b$$

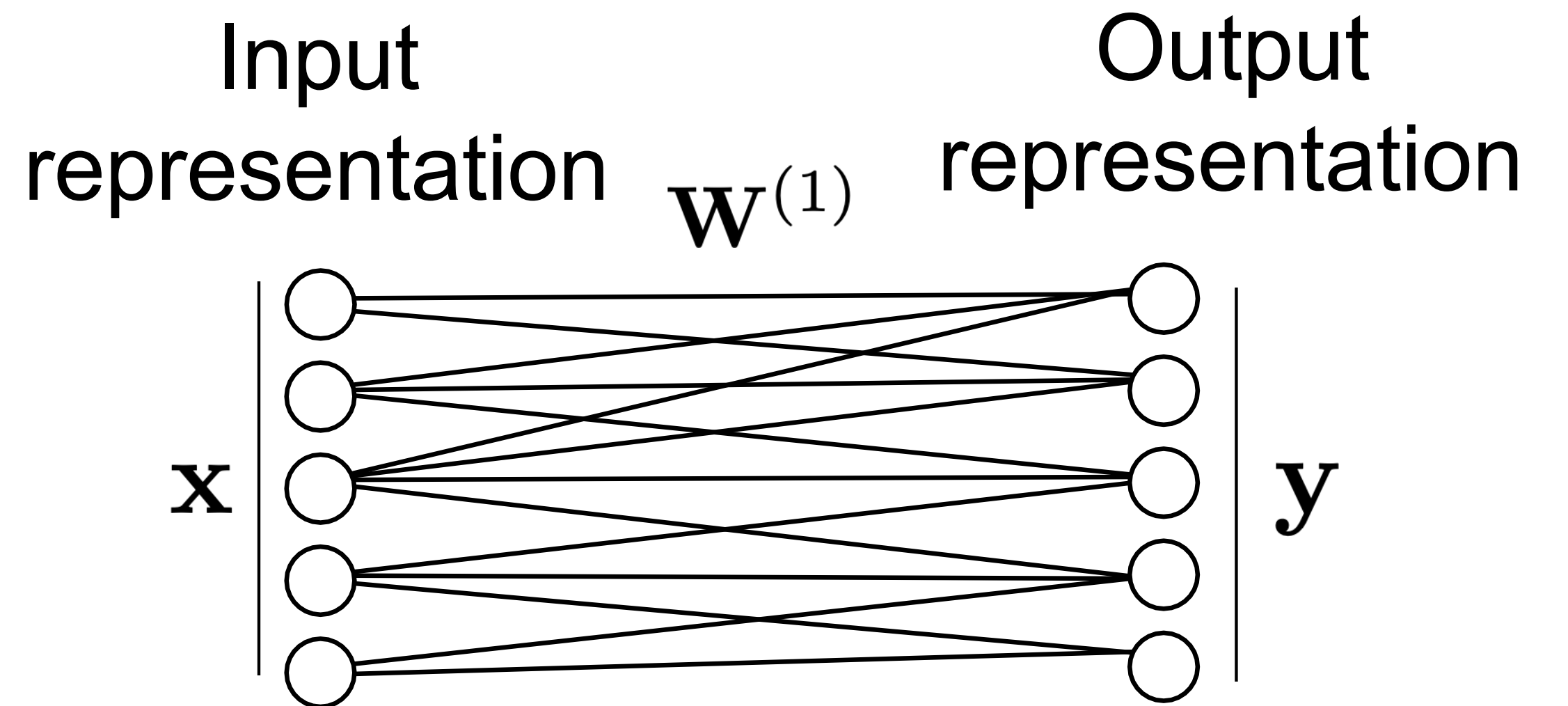
$$\begin{bmatrix} w_{11} & \cdots & w_{jn} & b_1 \\ \vdots & \ddots & \vdots & \vdots \\ w_{j1} & \cdots & w_{jn} & b_j \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_n \\ 1 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_j \end{bmatrix}$$

Can again simplify notation by appending a 1 to \mathbf{x}

Connectivity patterns



Fully connected layer

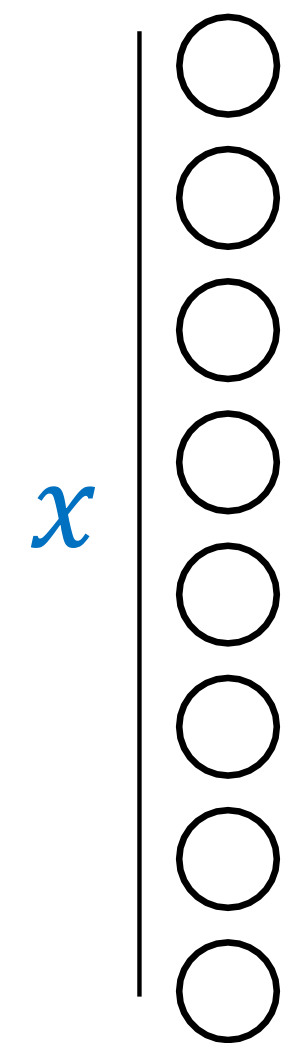


*Locally connected layer
(Sparse W)*

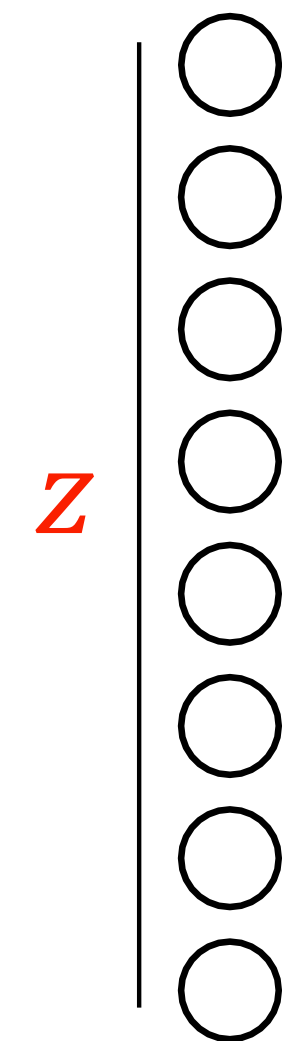
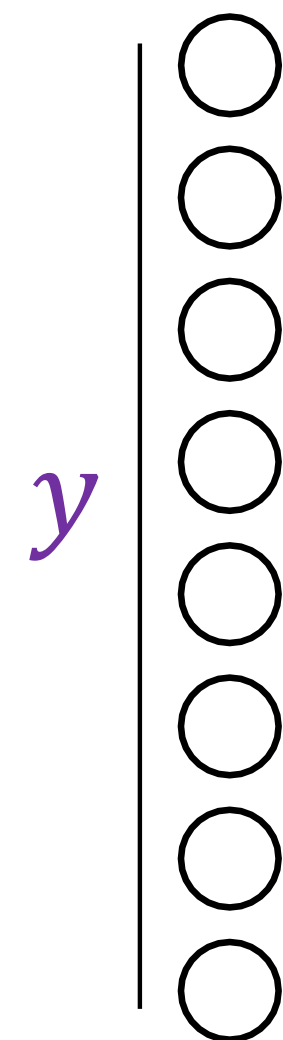
Computation in a neural net – Recap

We can now transform our input representation vector into some output representation vector using a bunch of linear combinations of the input:

Input
representation

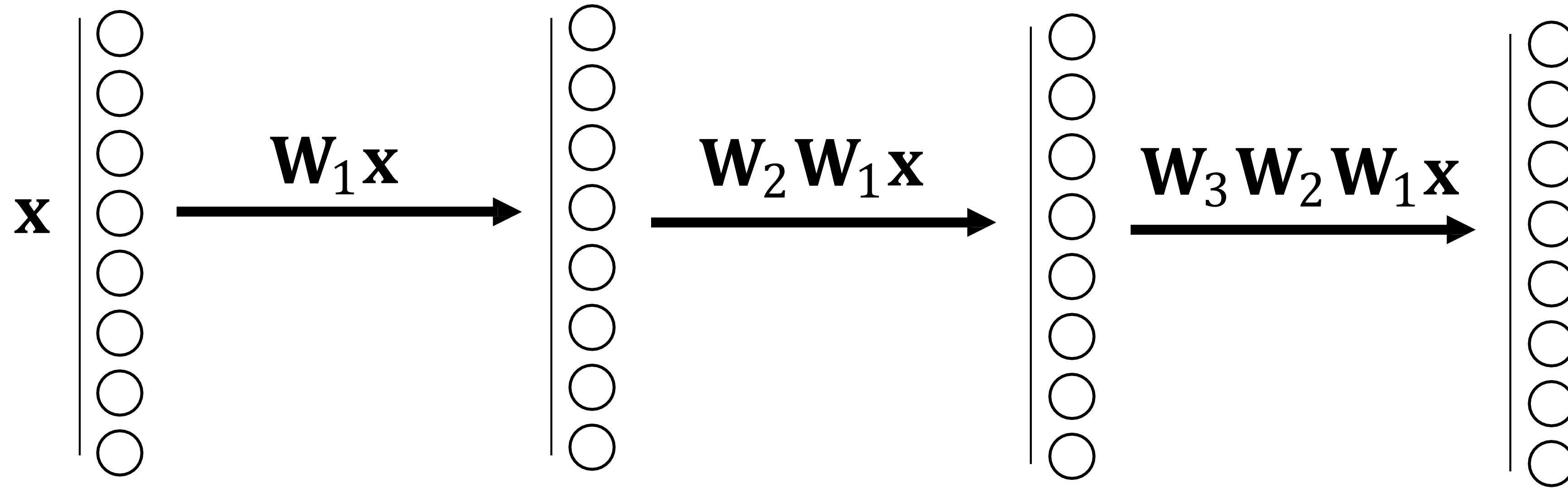


Output
representation

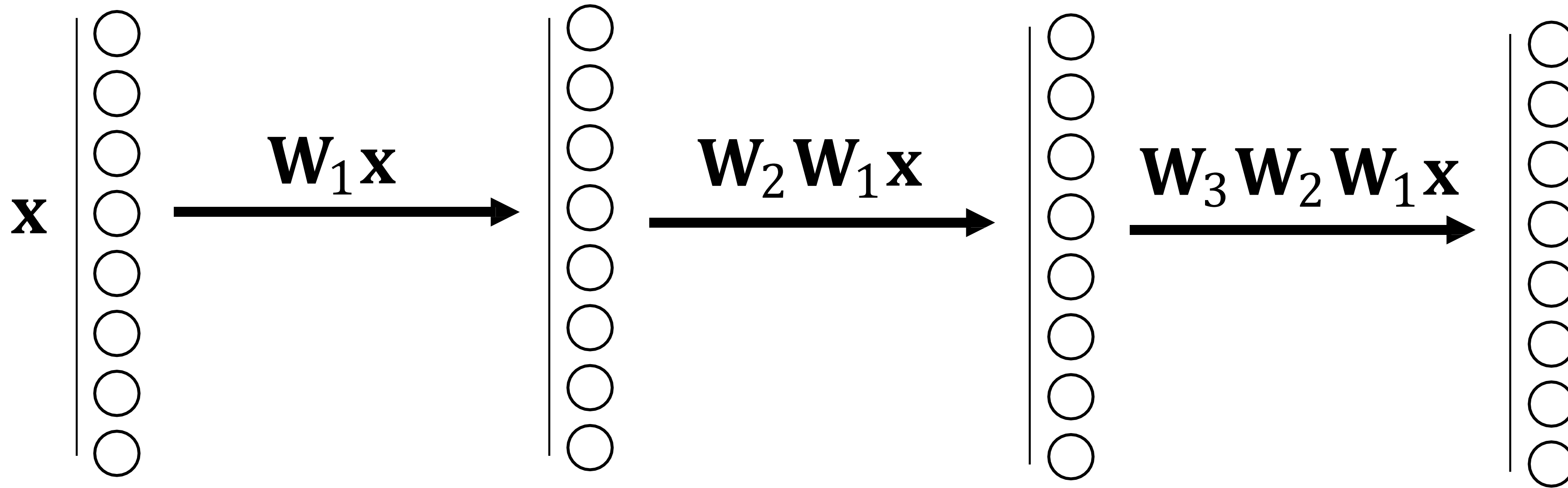


We can repeat this as
many times as we want!

What is the problem with this idea?



What is the problem with this idea?



Can be expressed as single linear layer!

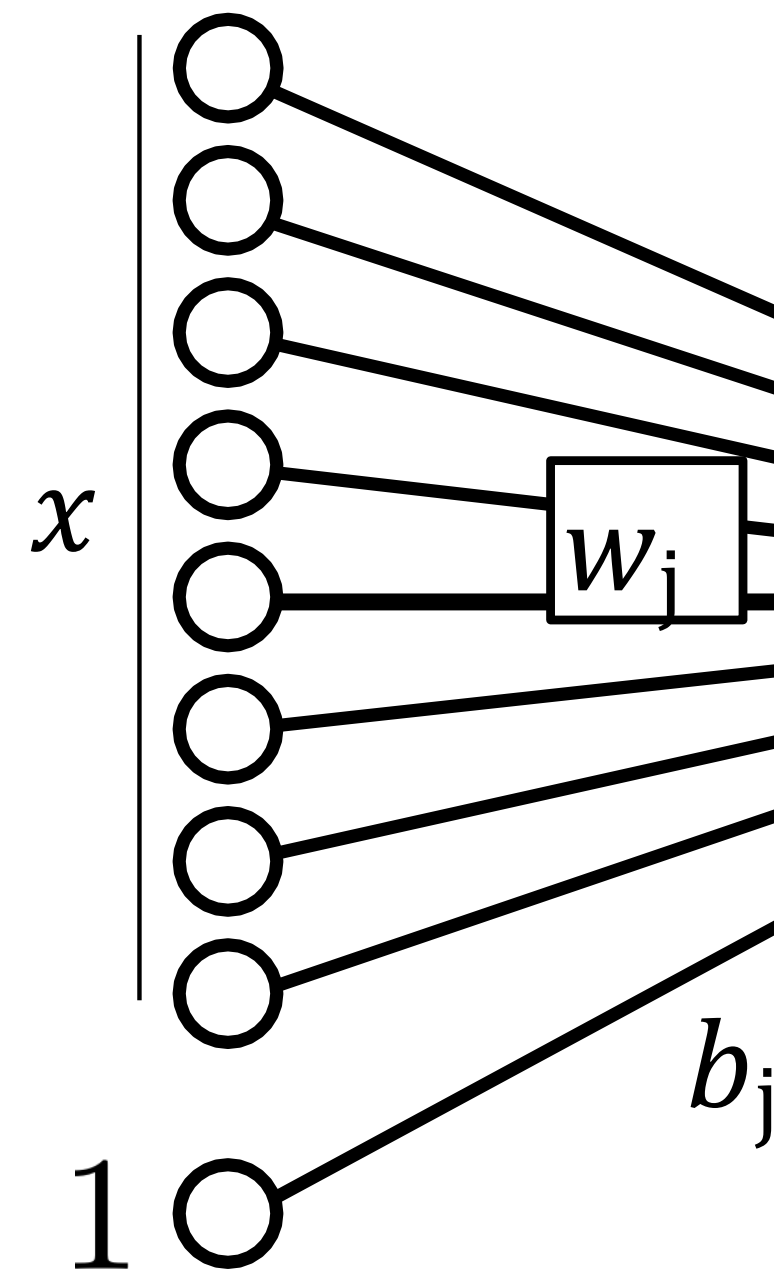
$$\widehat{W}x$$

Limited power: can't solve XOR ☹️

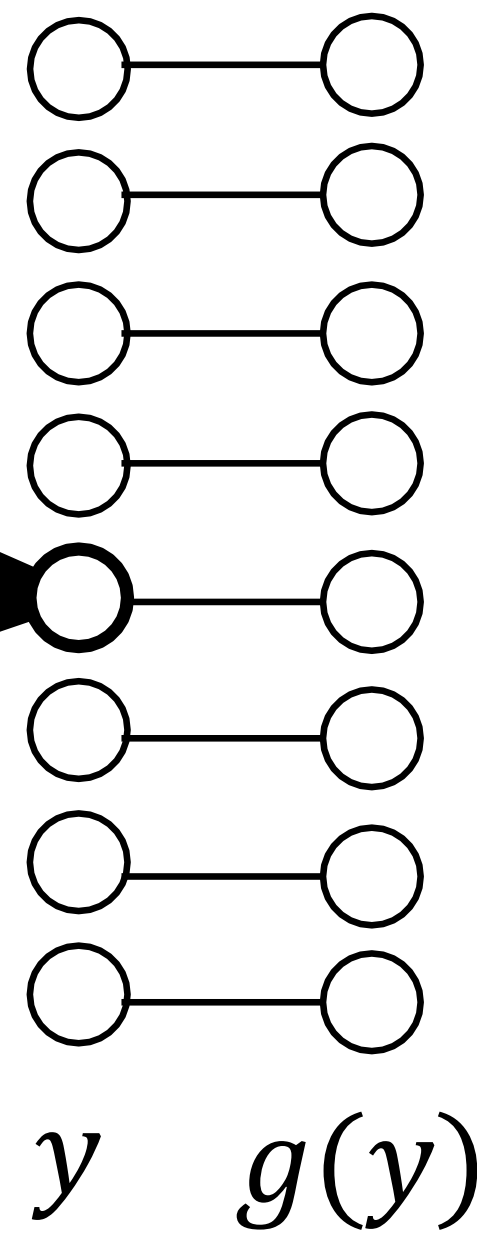
Solution: simple nonlinearity

Linear layer

Input representation



Output representation

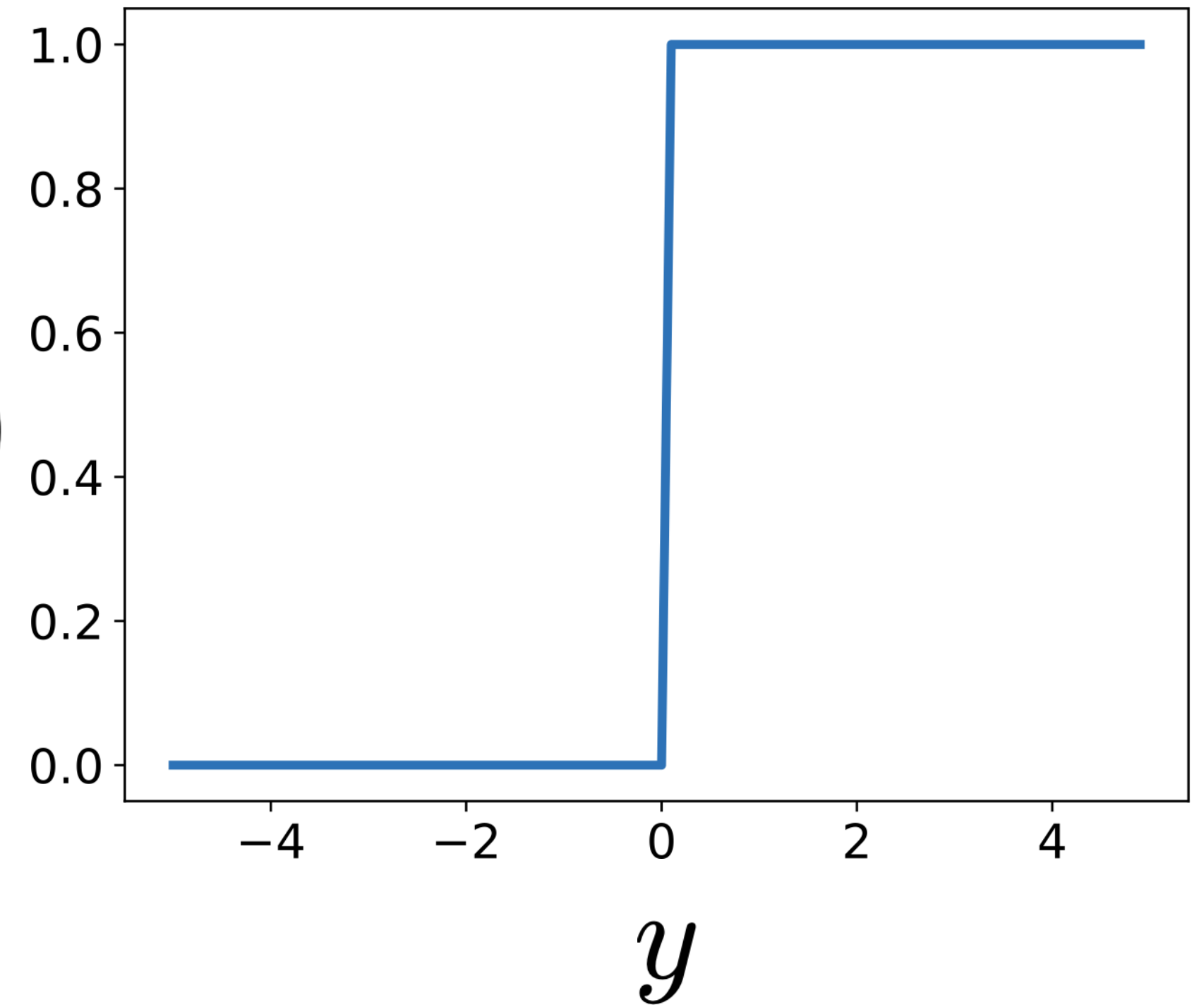


w_j

b_j

$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$

$g(y)$



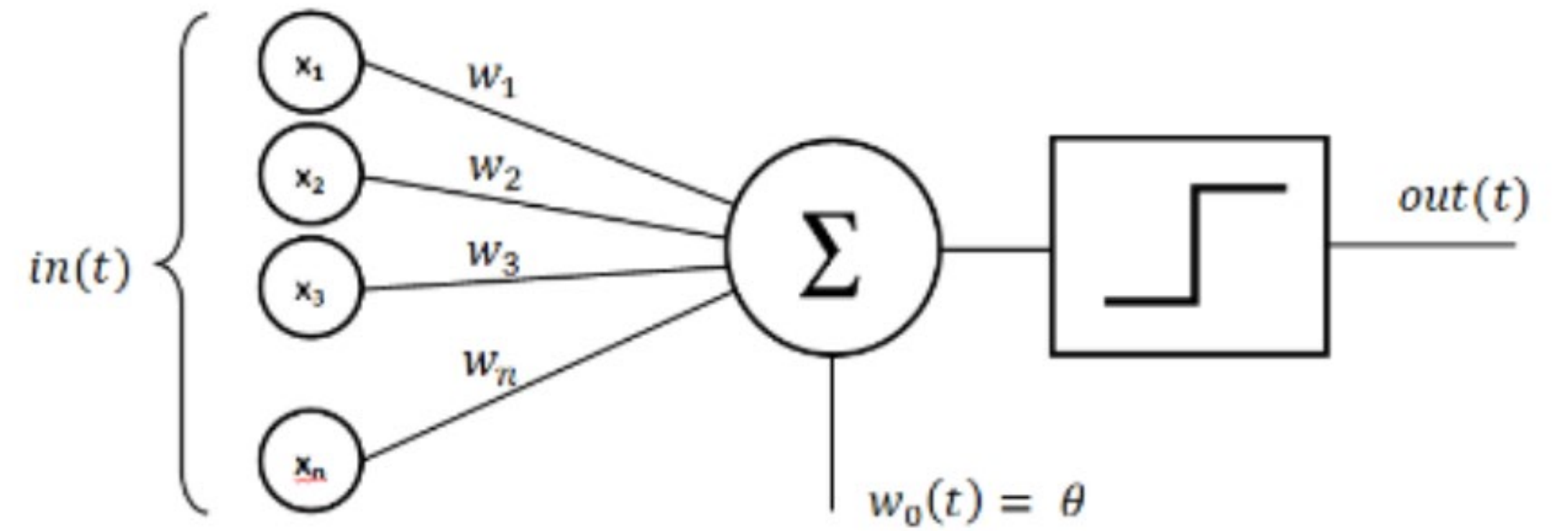
Pointwise
Non-linearity

The Perceptron

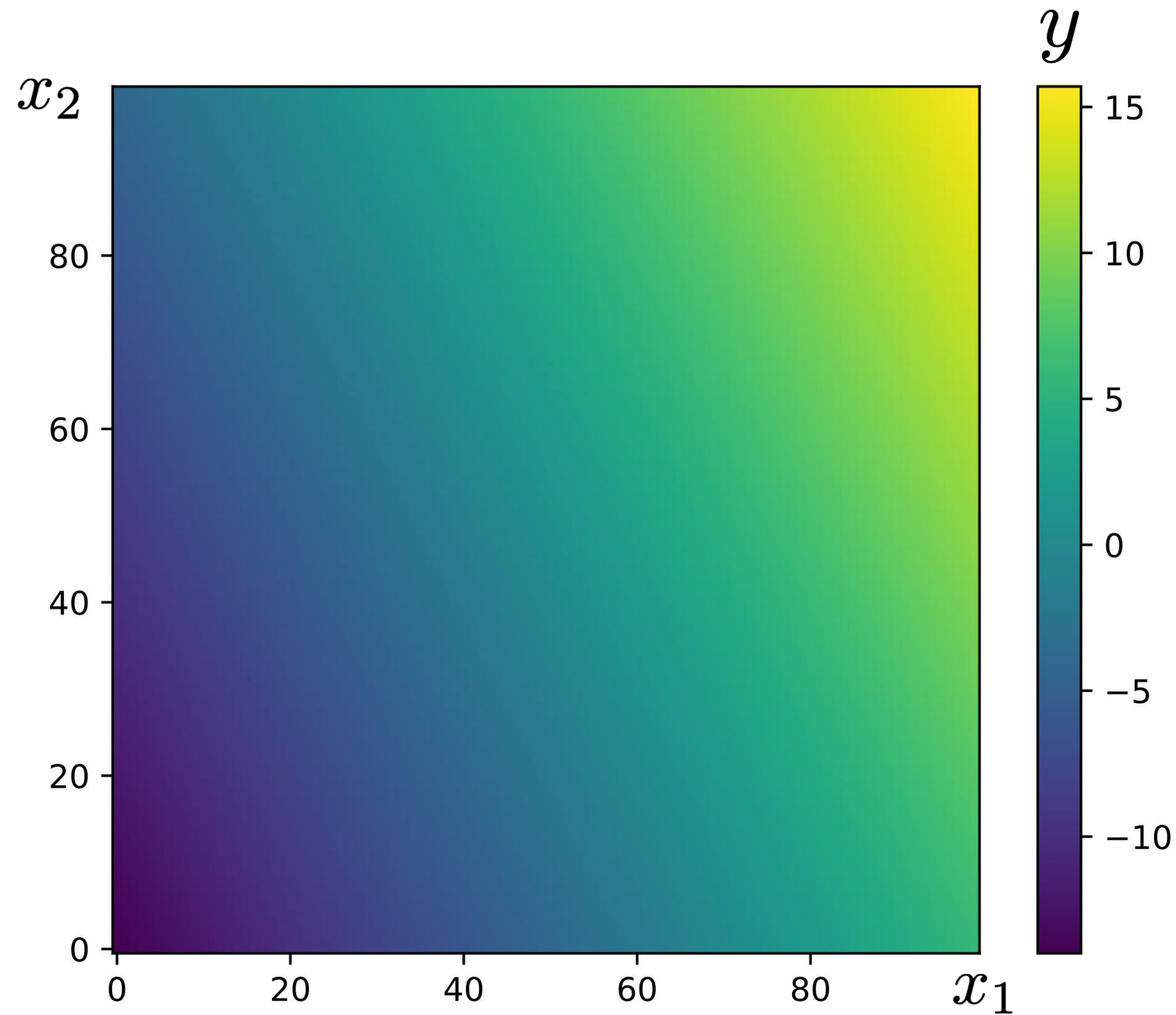
WHAT PERCEPTRON SOUNDS LIKE



wHat pERceptRoNS ArE

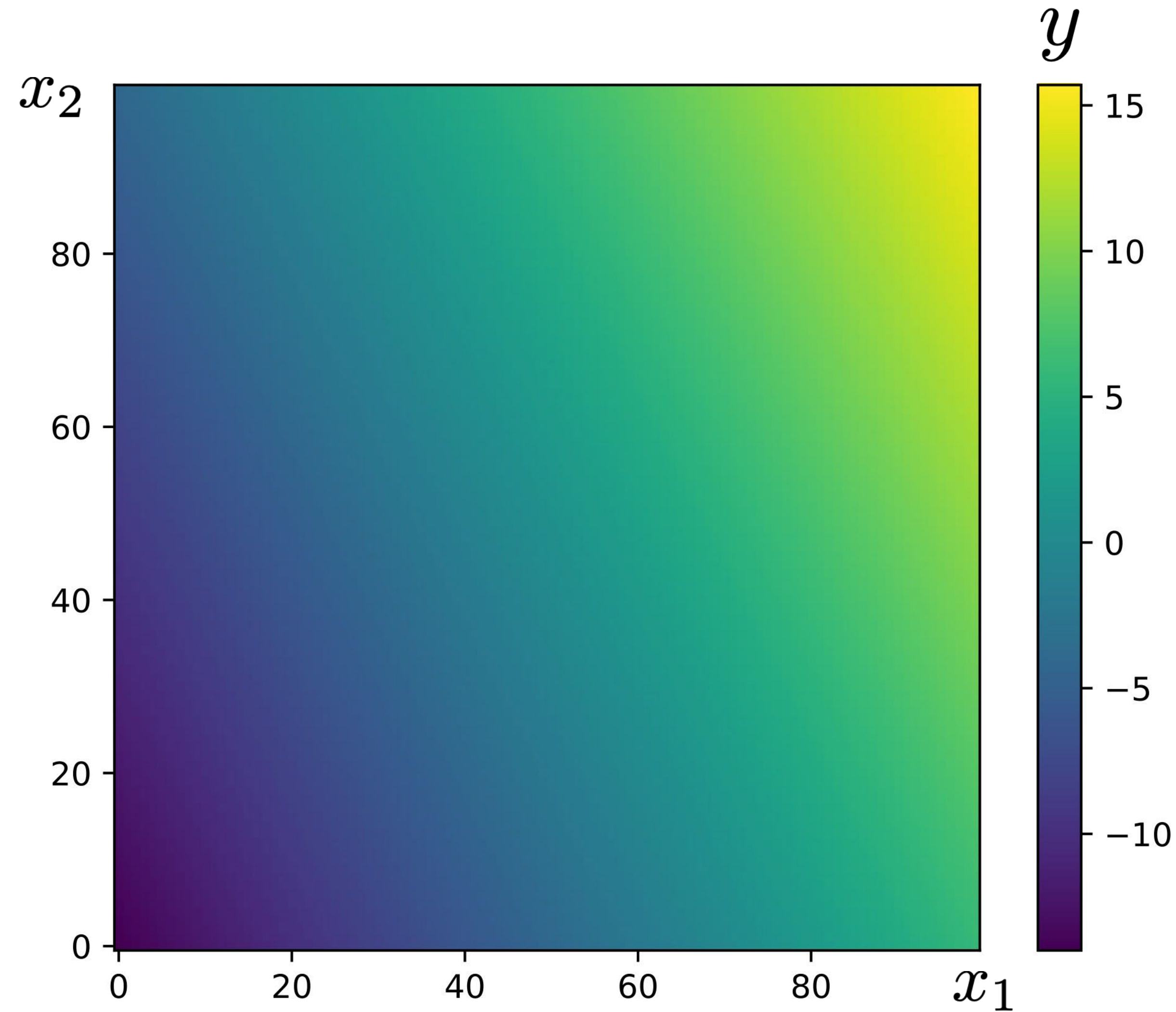


Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

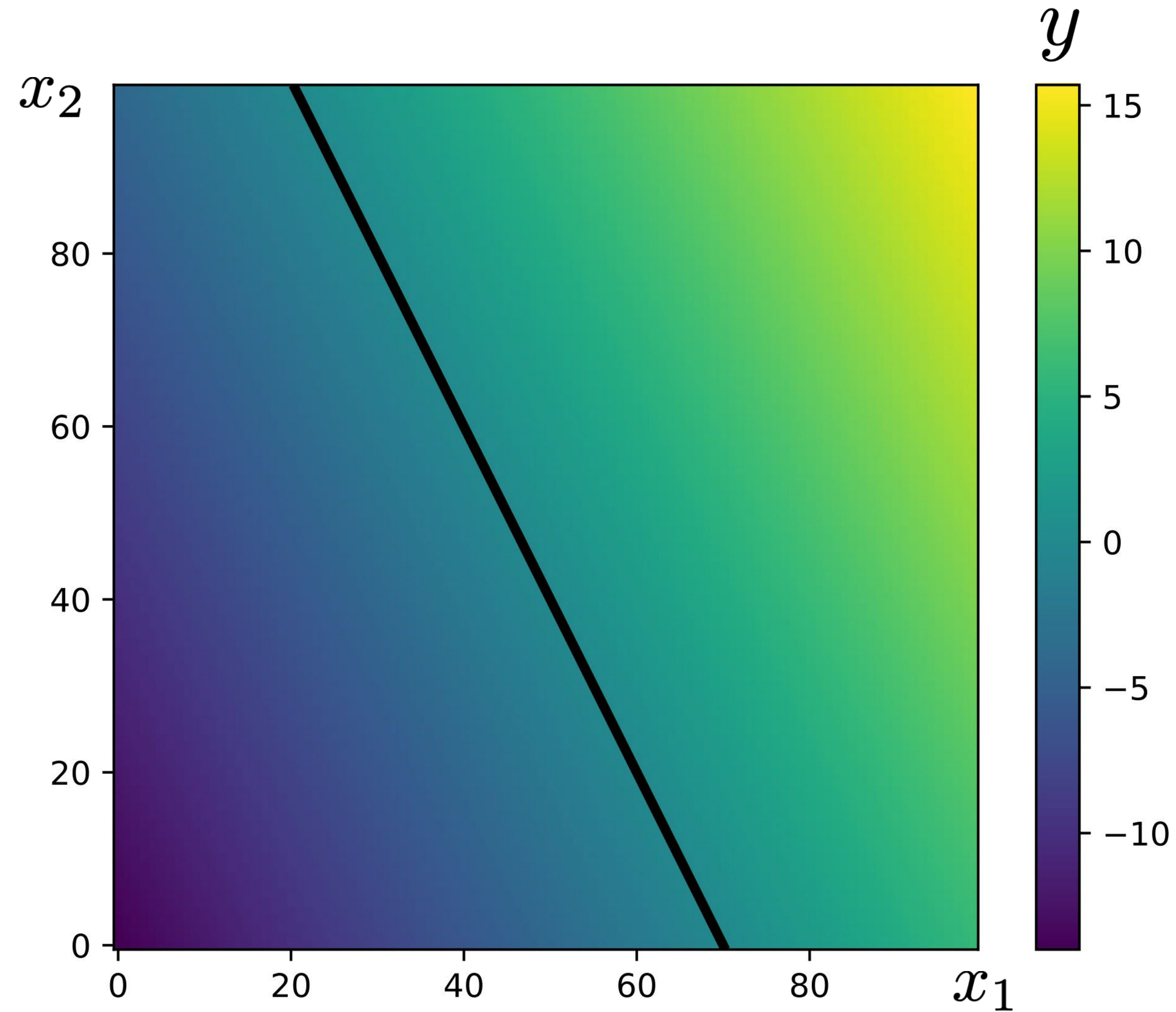
Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$

Example: linear classification with a perceptron



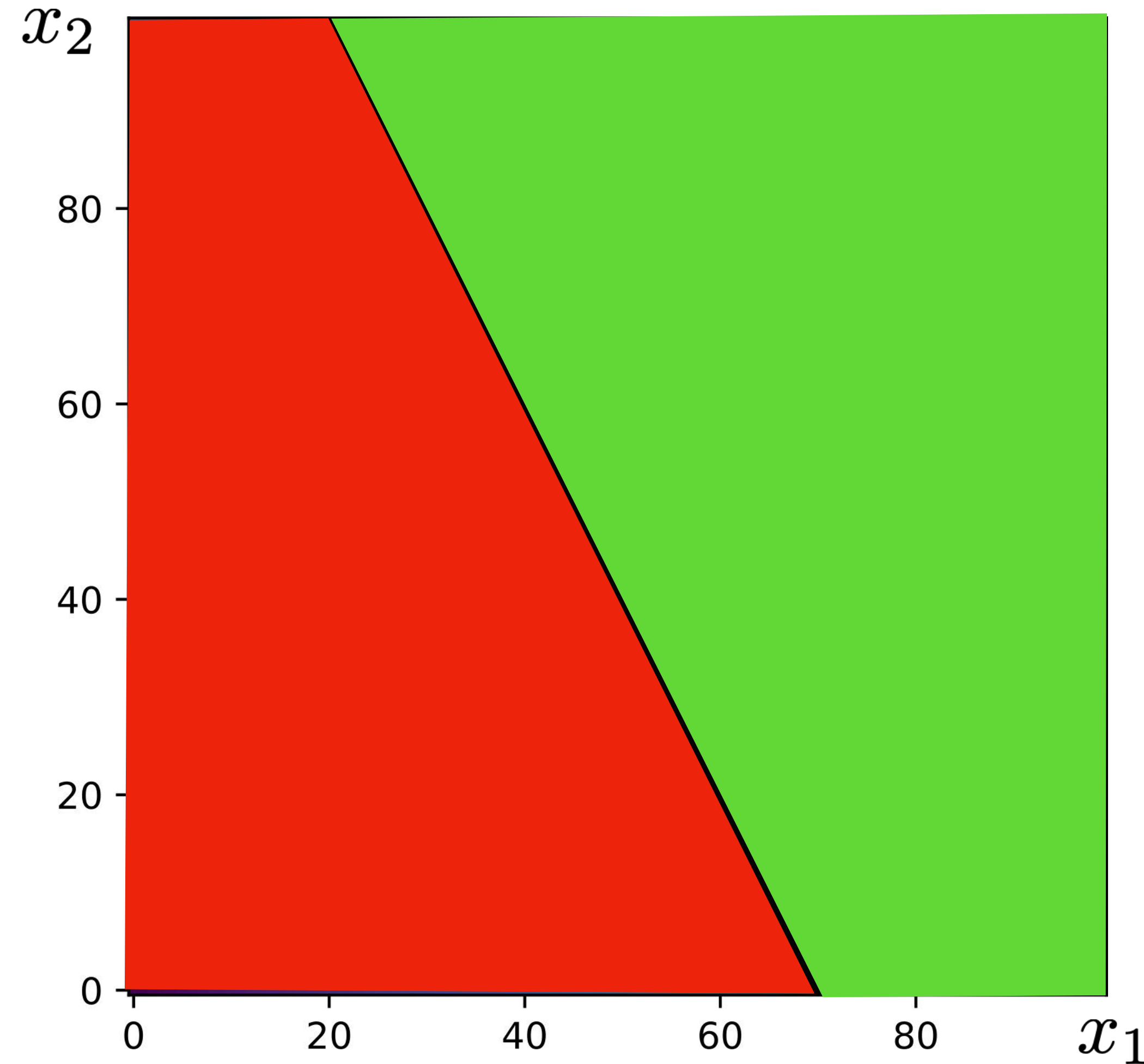
$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$

“when y is greater than 0, set all pixel values to 1 (green), otherwise, set all pixel values to 0 (red)”

Example: linear classification with a perceptron

$$g(y)$$



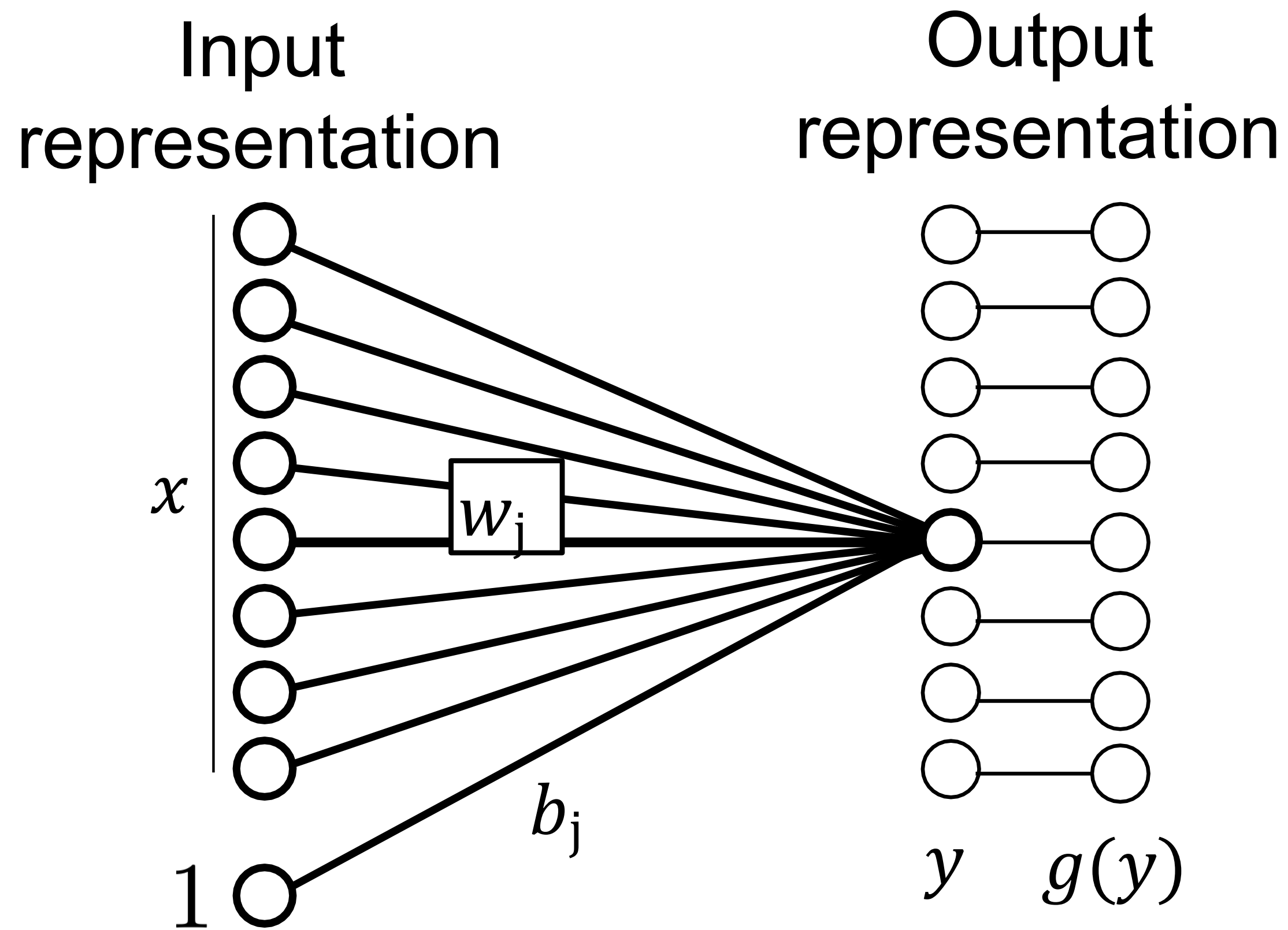
$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$

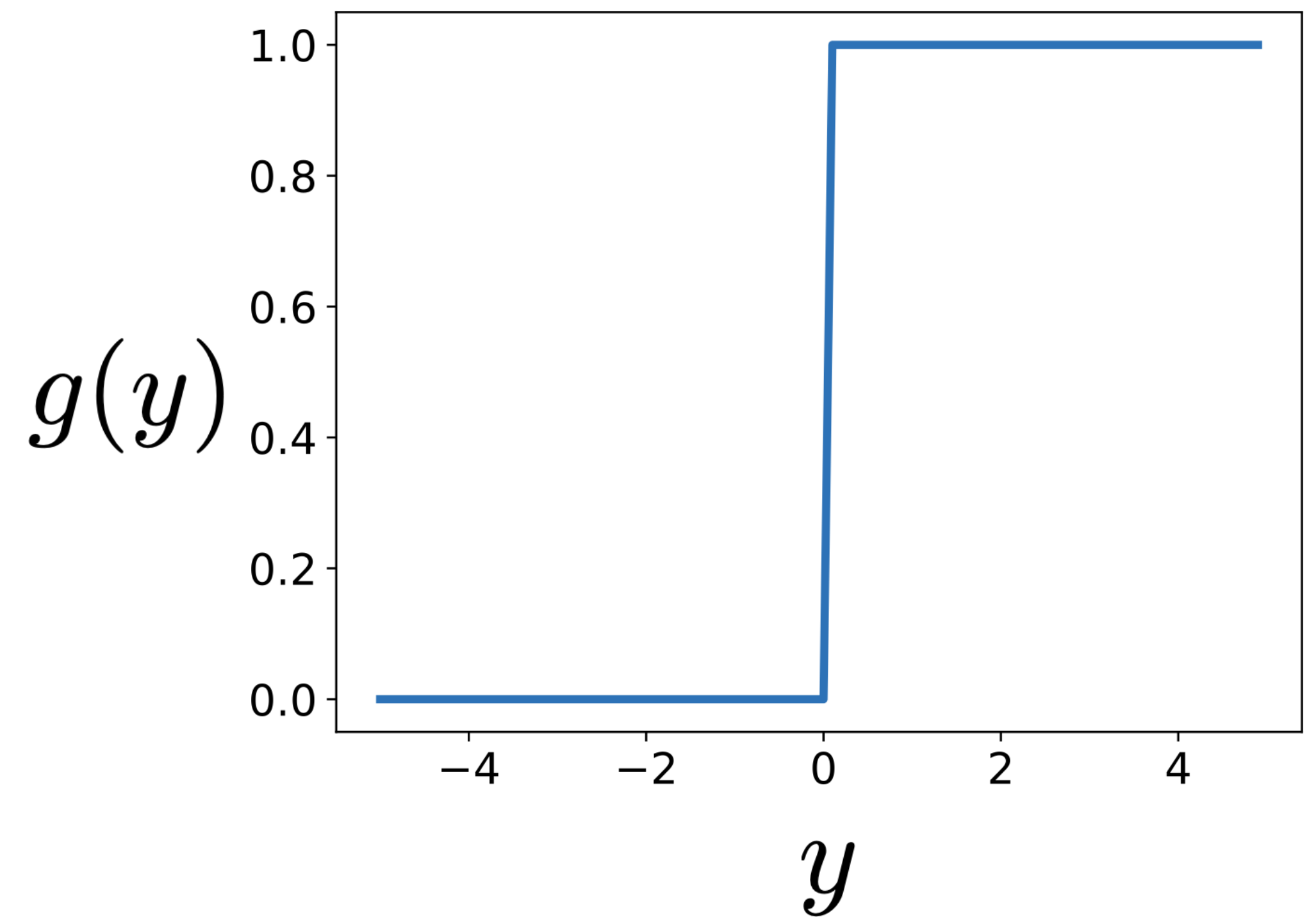
“when y is greater than 0, set all pixel values to 1 (green), otherwise, set all pixel values to 0 (red)”

Computation in a neural net - nonlinearity

Linear layer



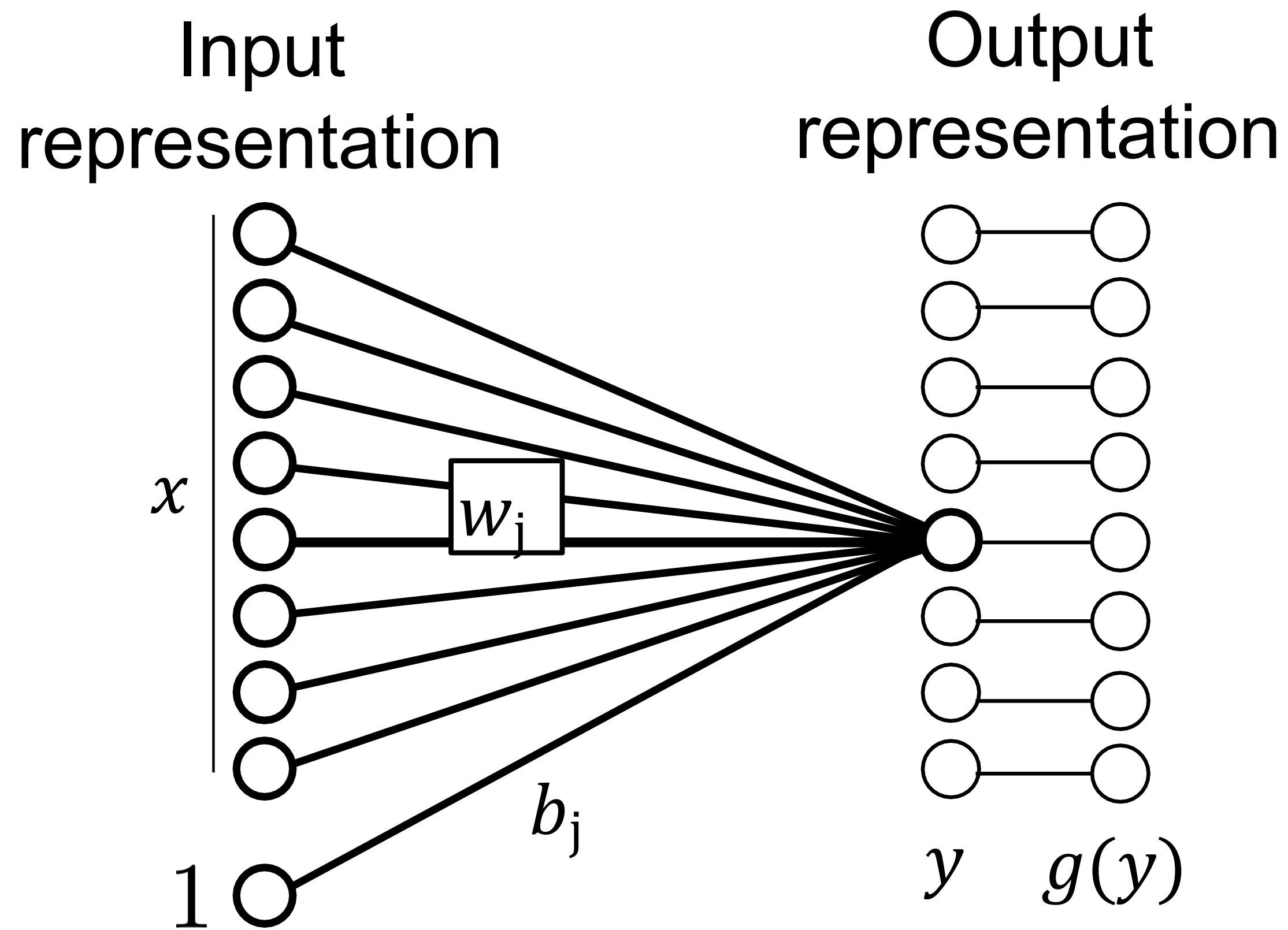
$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$



Can't use with gradient descent, $\frac{\partial}{\partial y} g = 0$

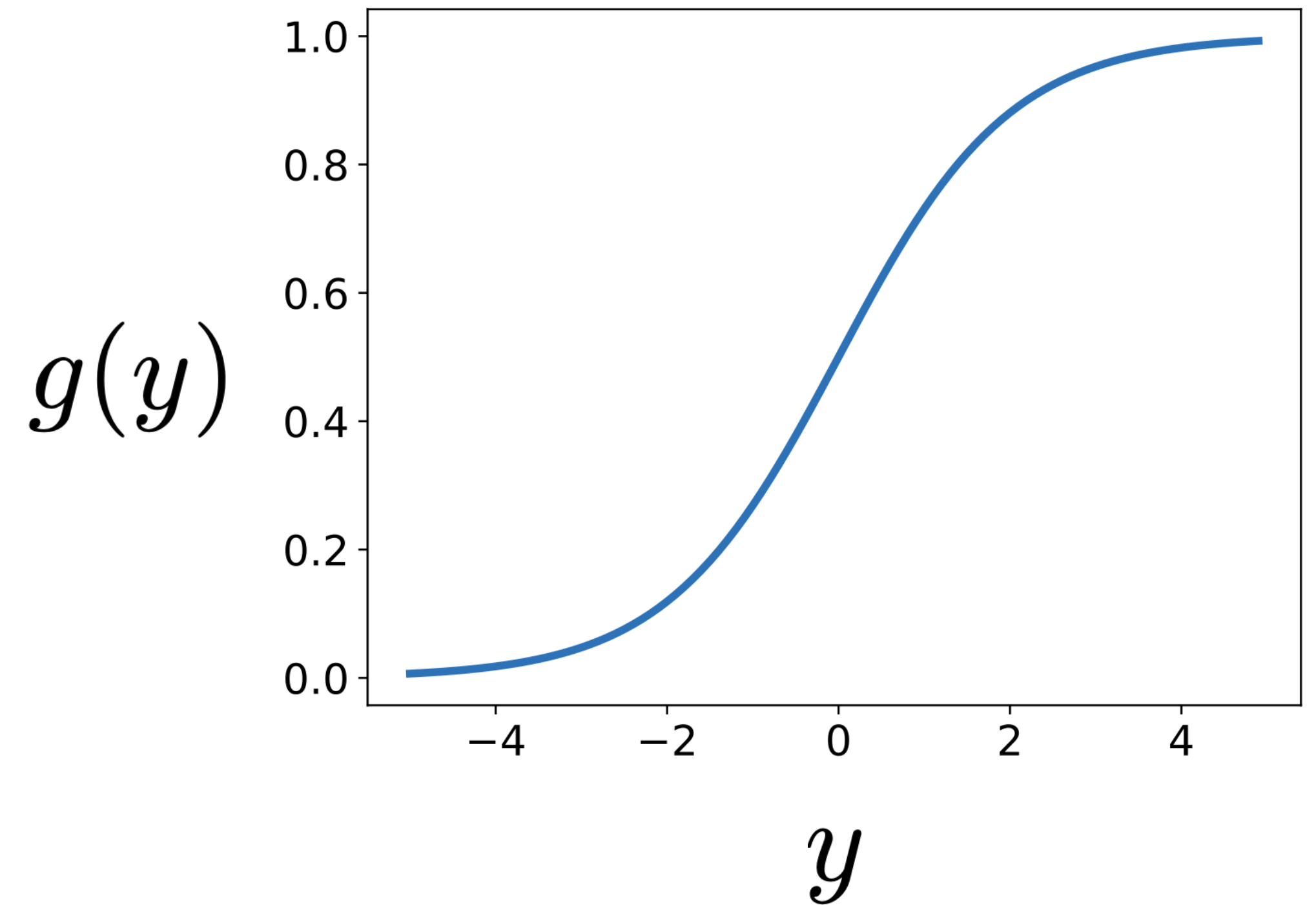
Computation in a neural net - nonlinearity

Linear layer



Sigmoid

$$g(y) = \sigma(y) = \frac{1}{1 + e^{-y}}$$

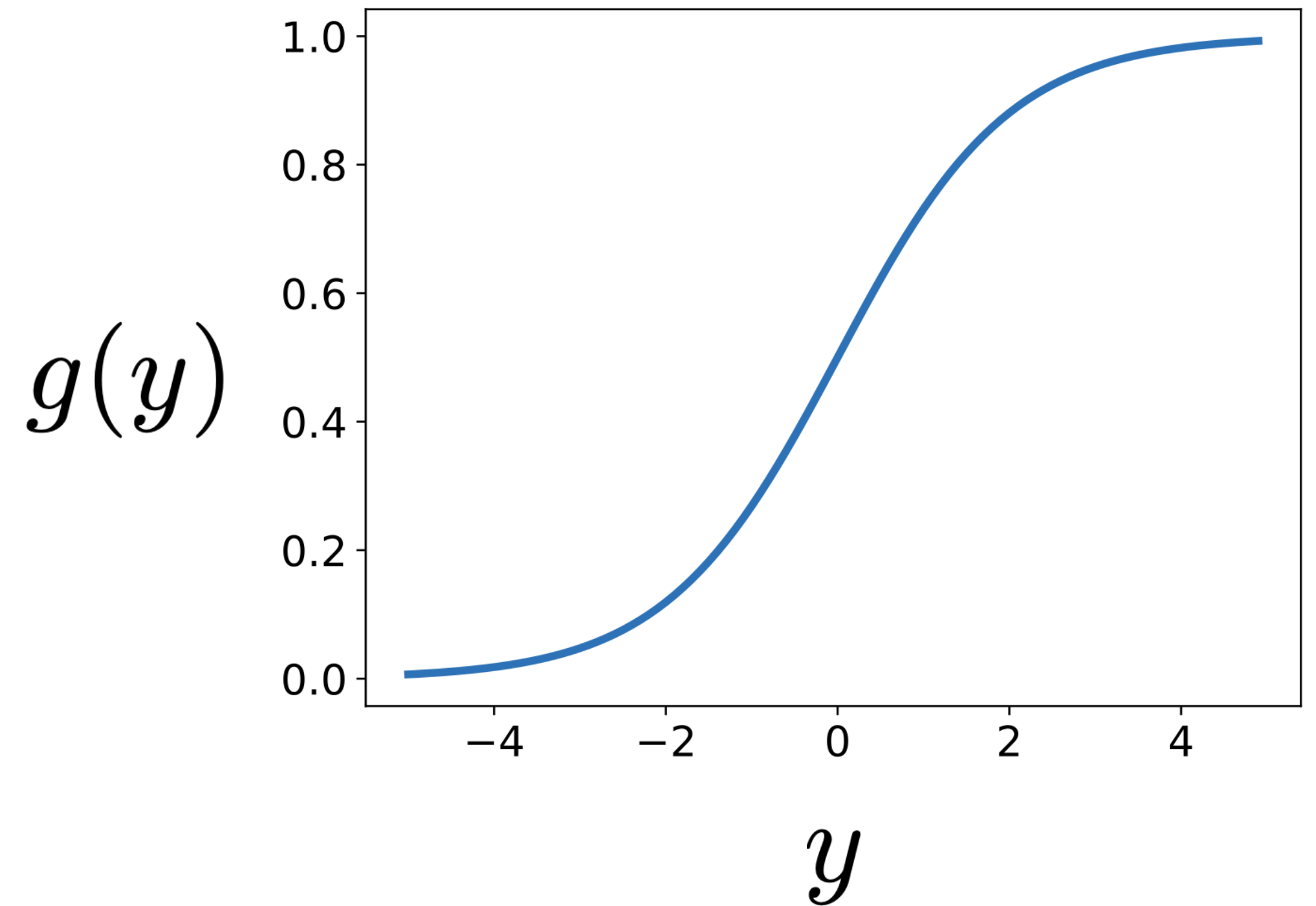


Computation in a neural net - nonlinearity

- Bounded between [0,1]
- Saturation for large +/- inputs
- Gradients go to zero

Sigmoid

$$g(y) = \sigma(y) = \frac{1}{1 + e^{-y}}$$



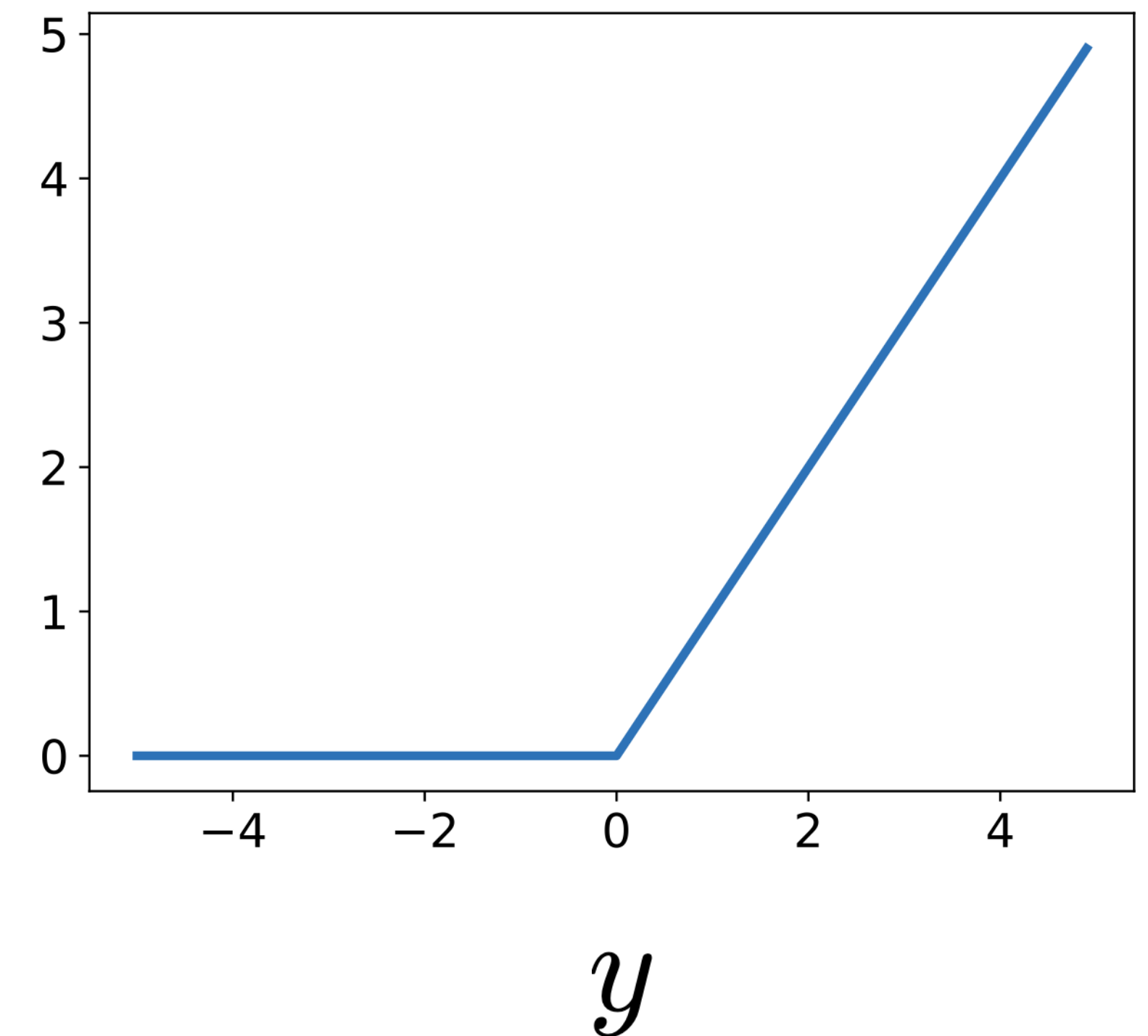
Computation in a neural net — nonlinearity

- Unbounded output (on positive side)
- Efficient to implement: $\frac{\partial g}{\partial y} = \begin{cases} 0, & \text{if } y < 0 \\ 1, & \text{if } y \geq 0 \end{cases}$
- Also seems to help convergence (6x speedup vs. tanh in [Krizhevsky et al. 2012])
- Drawback: if strongly in negative region, unit is dead forever (no gradient).
- Default choice: widely used in current models!

Rectified linear unit (ReLU)

$$g(y) = \max(0, y)$$

$g(y)$



Computation in a neural net — nonlinearity

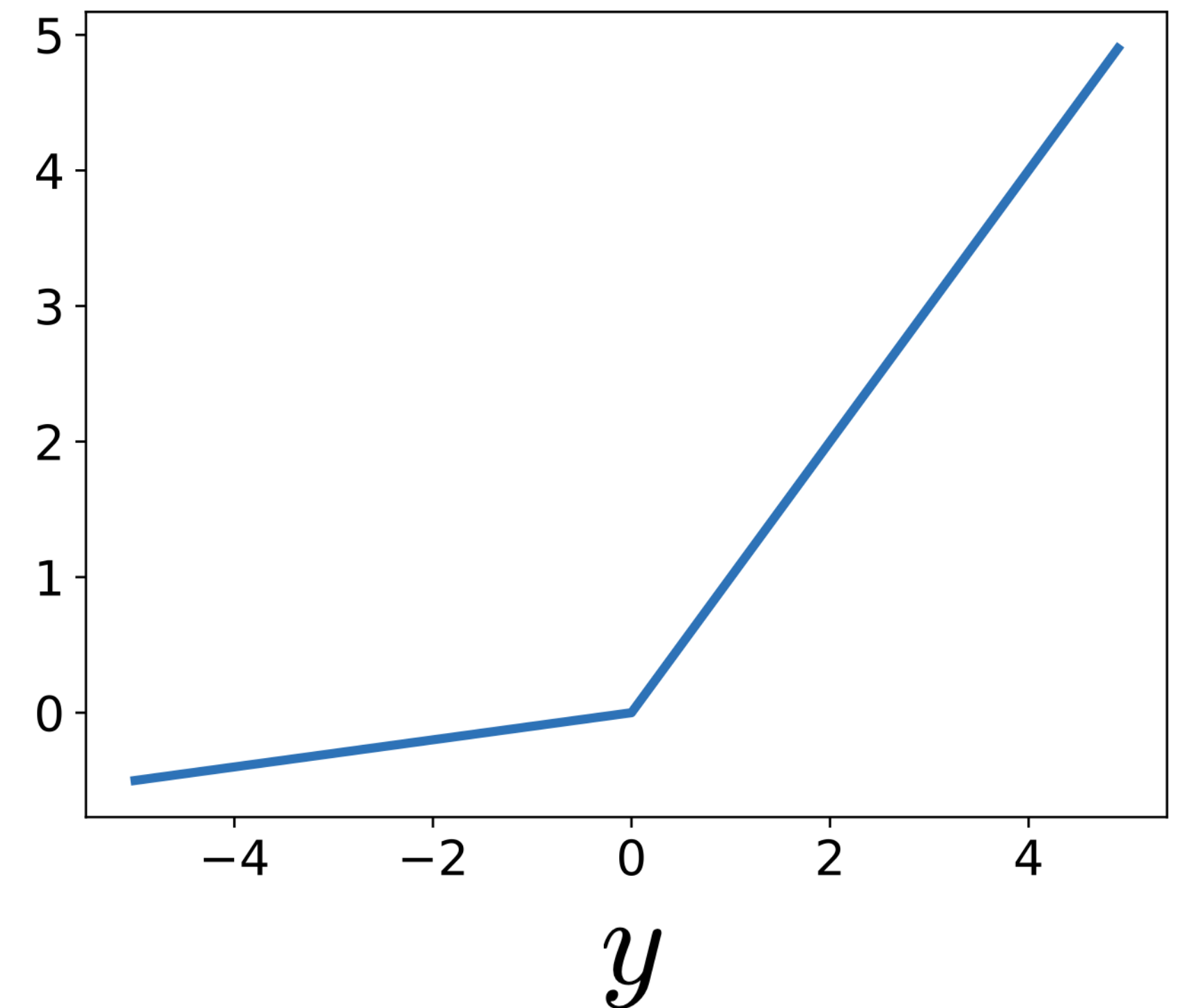
- where a is small (e.g., 0.02)
- Efficient to implement:
- Has non-zero gradients everywhere (unlike ReLU)

$$\frac{\partial g}{\partial y} = \begin{cases} -a, & \text{if } y < 0 \\ 1, & \text{if } y \geq 0 \end{cases}$$

Leaky ReLU

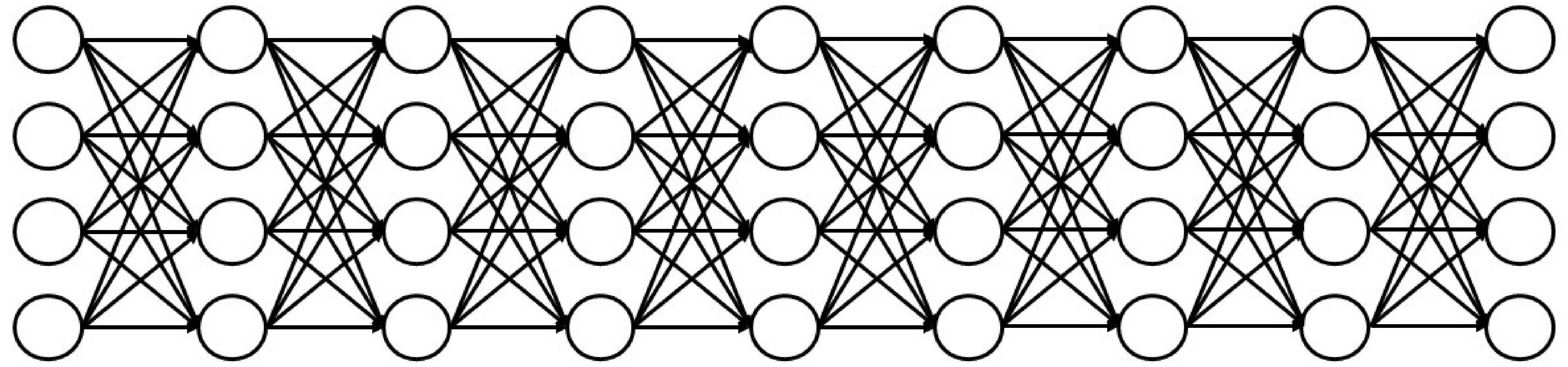
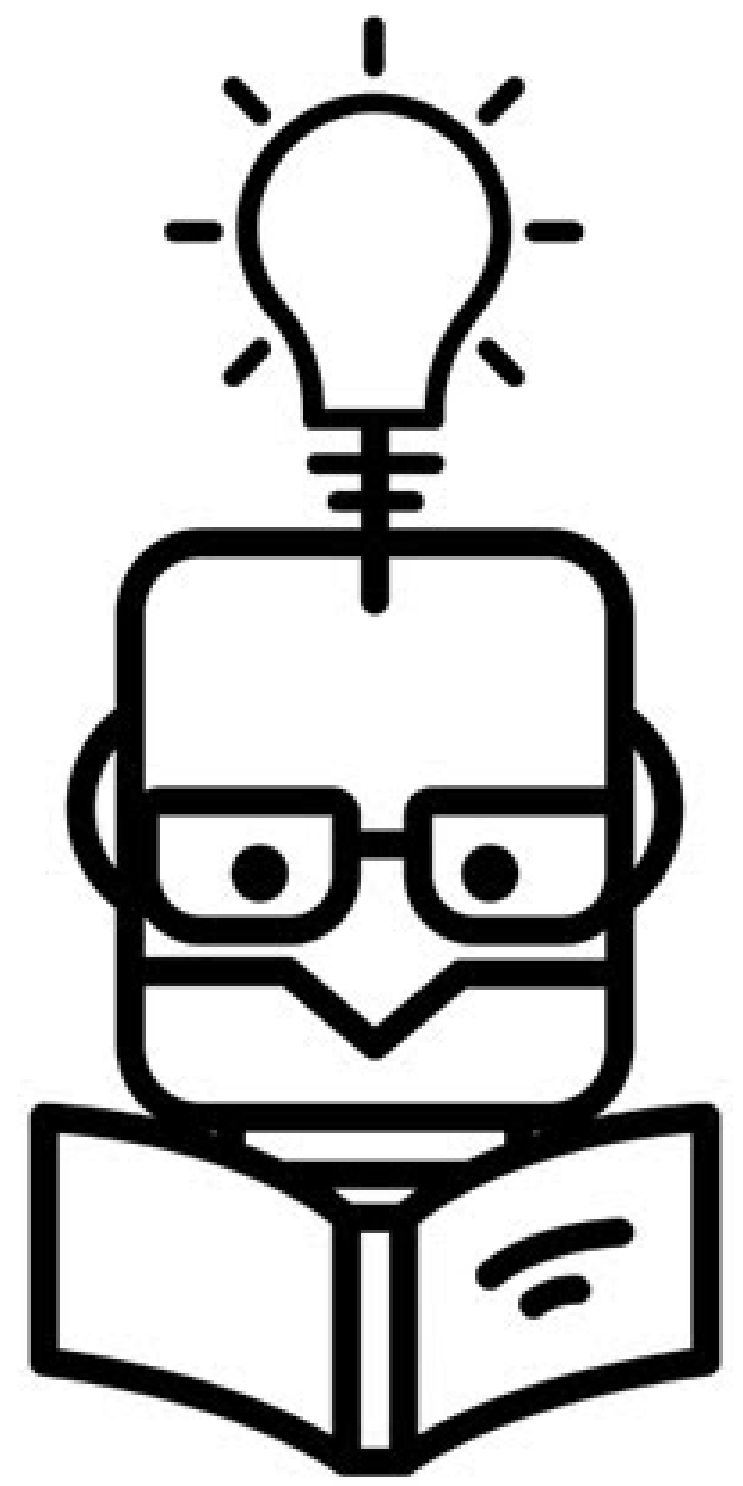
$$g(y) = \begin{cases} \max(0, y), & \text{if } y \geq 0 \\ a \min(0, y), & \text{if } y < 0 \end{cases}$$

$g(y)$





DEEP Neural Nets?

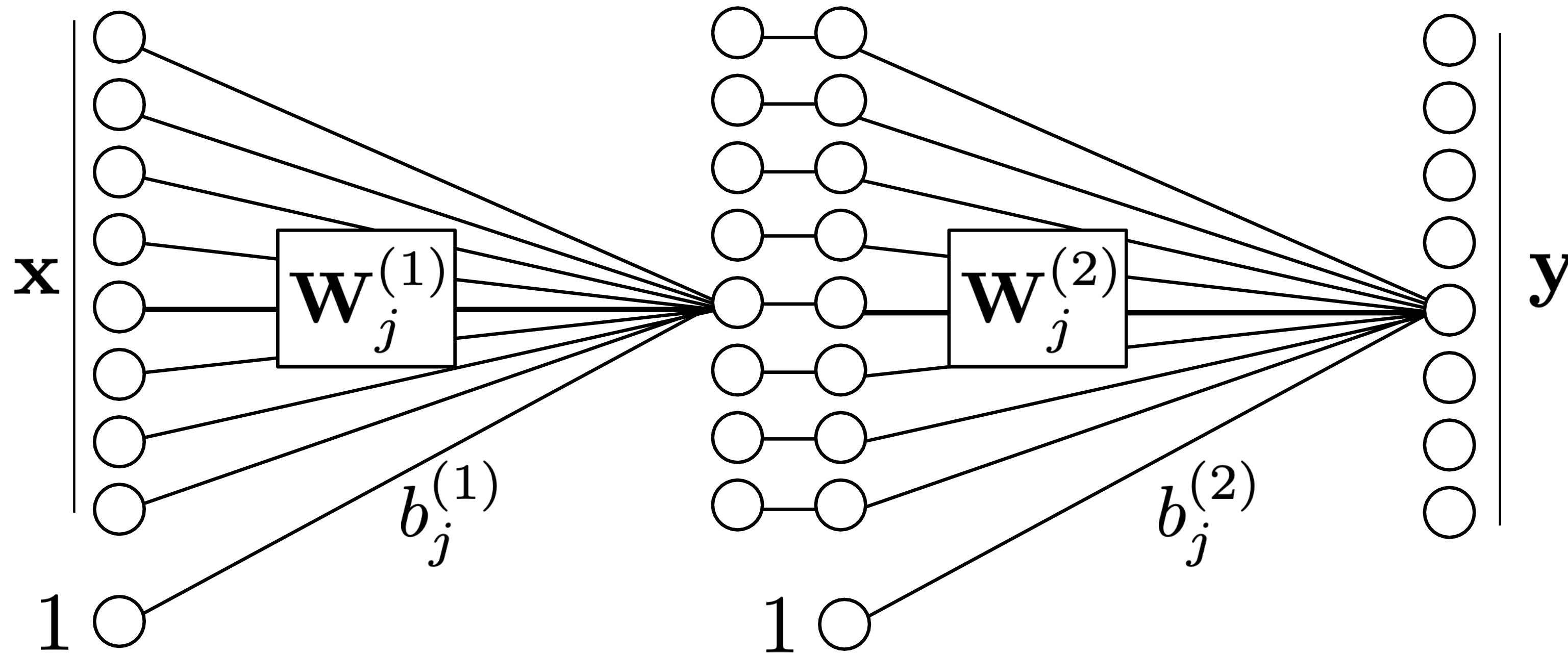


Stacking layers

Input
representation

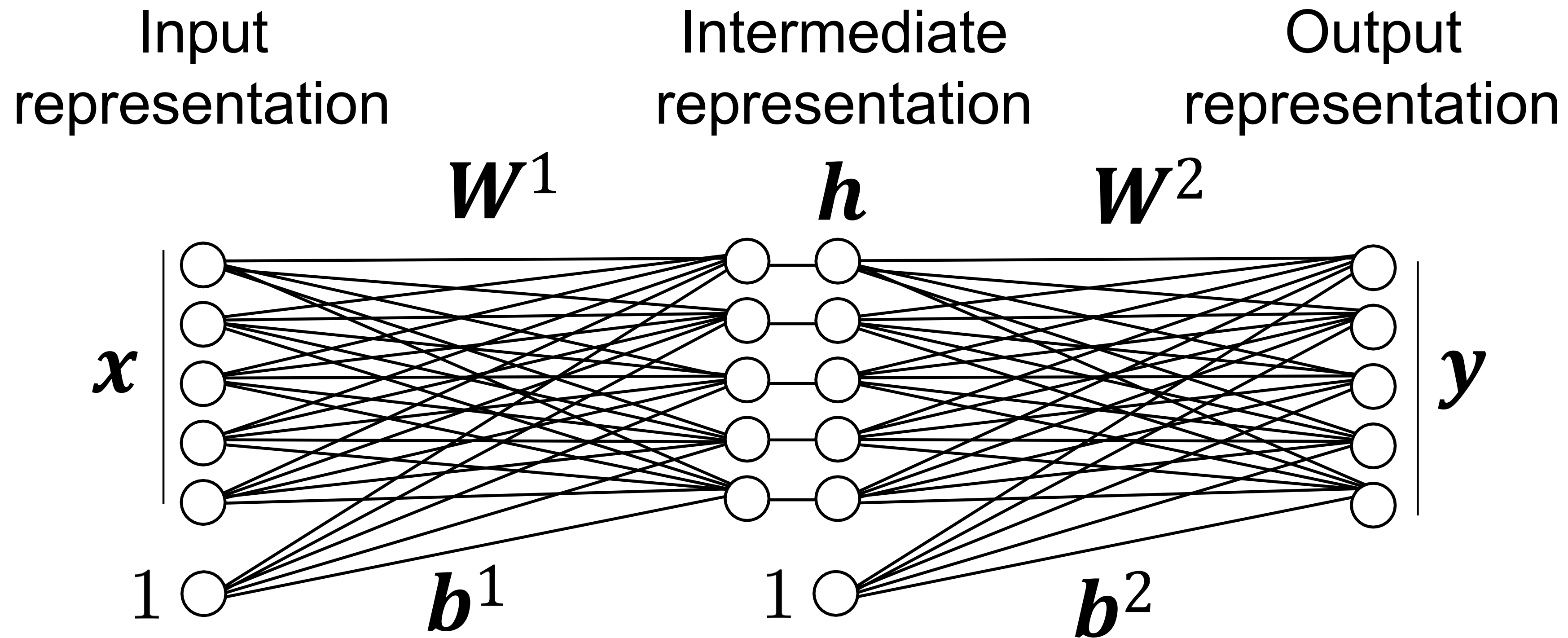
Intermediate
representation

Output
representation



\mathbf{h} = "hidden units"

Stacking layers

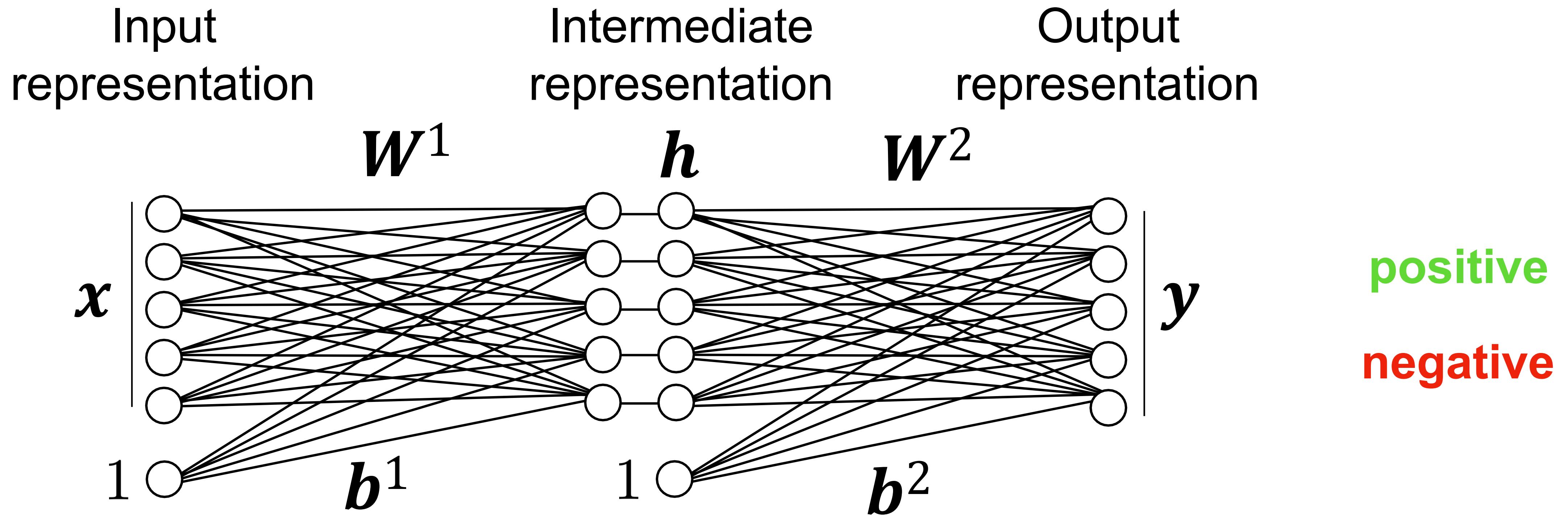


$$h = g(W^1 x + b^1) \quad y = g(W^2 h + b^2)$$

ReLU \nearrow

$$\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$$

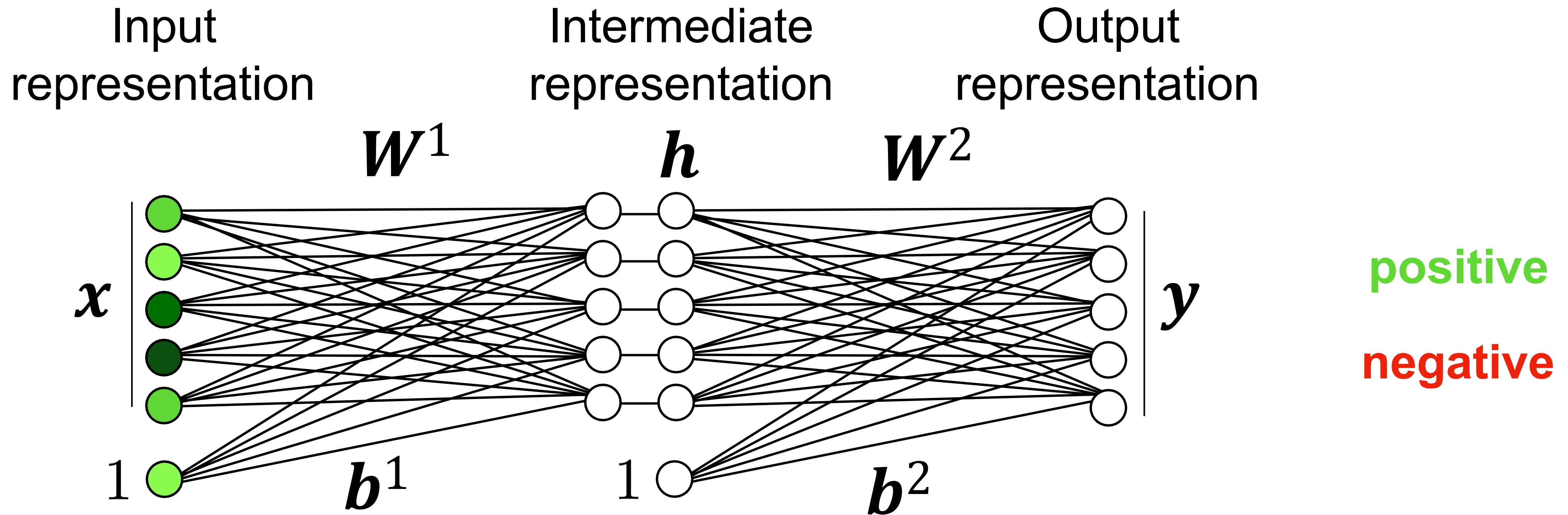
Stacking layers



$$h = g(W^1 x + b^1) \quad y = g(W^2 h + b^2)$$

ReLU \nearrow $\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$

Stacking layers

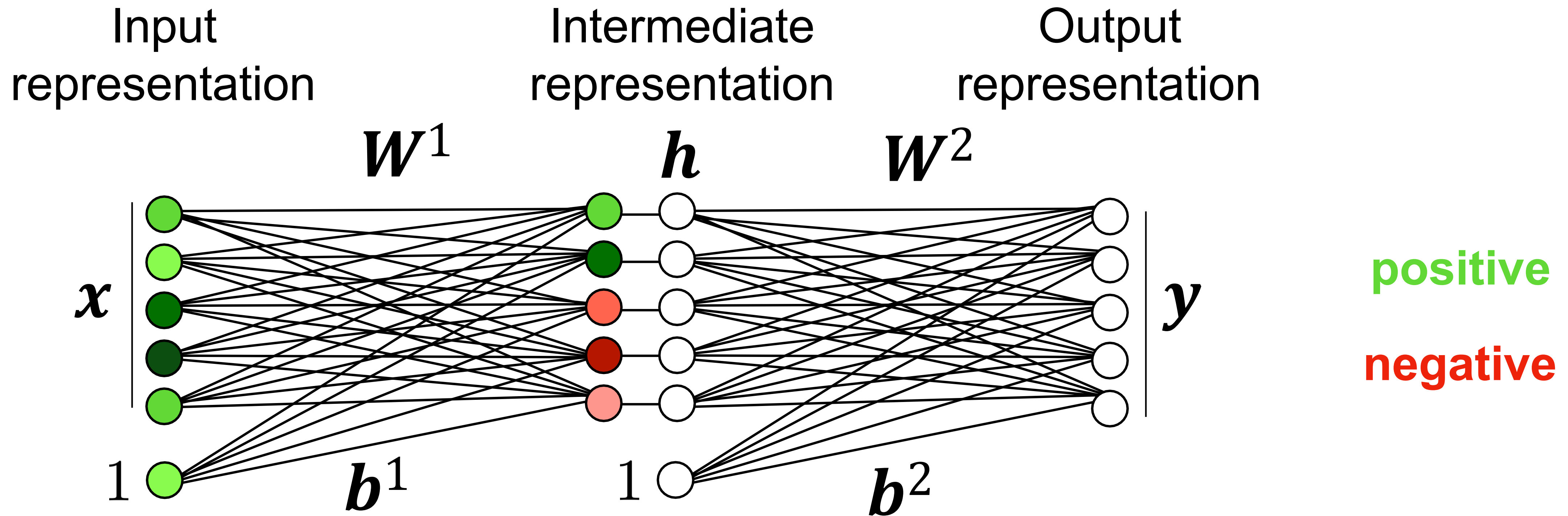


$$h = g(W^1 x + b^1) \quad y = g(W^2 h + b^2)$$

ReLU \nearrow

$$\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$$

Stacking layers



$$h = g(W^1 x + b^1) \quad y = g(W^2 h + b^2)$$

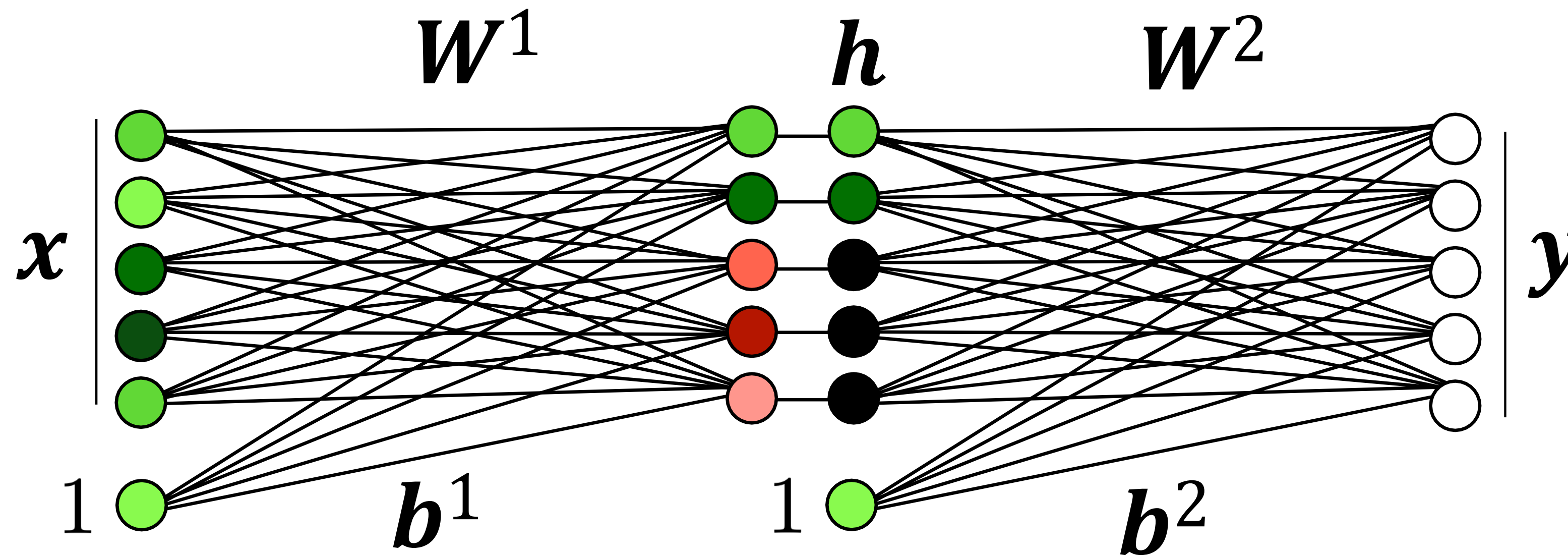
ReLU \nearrow $\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$

Stacking layers

Input representation

Intermediate representation

Output representation



positive

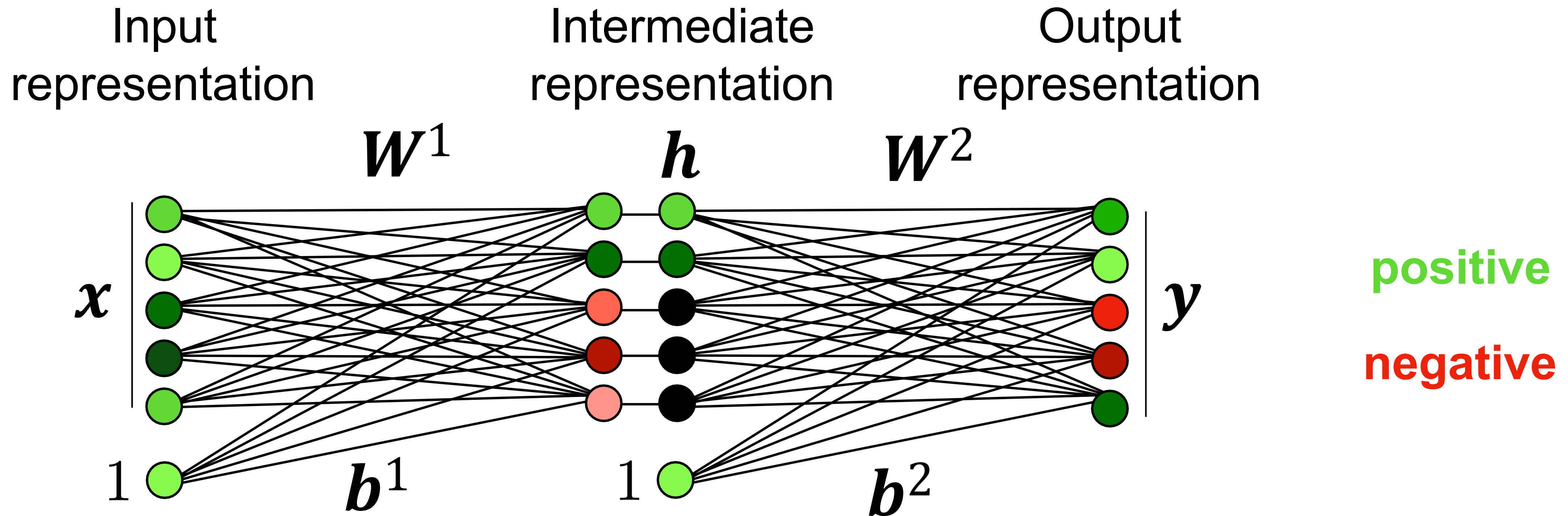
negative

$$\mathbf{h} = g(\mathbf{W}^1 \mathbf{x} + \mathbf{b}^1) \quad \mathbf{y} = g(\mathbf{W}^2 \mathbf{h} + \mathbf{b}^2)$$

ReLU

$$\theta = \{\mathbf{W}^1, \dots, \mathbf{W}^L, \mathbf{b}^1, \dots, \mathbf{b}^L\}$$

Stacking layers



$$h = g(W^1 x + b^1) \quad y = g(W^2 h + b^2)$$

ReLU

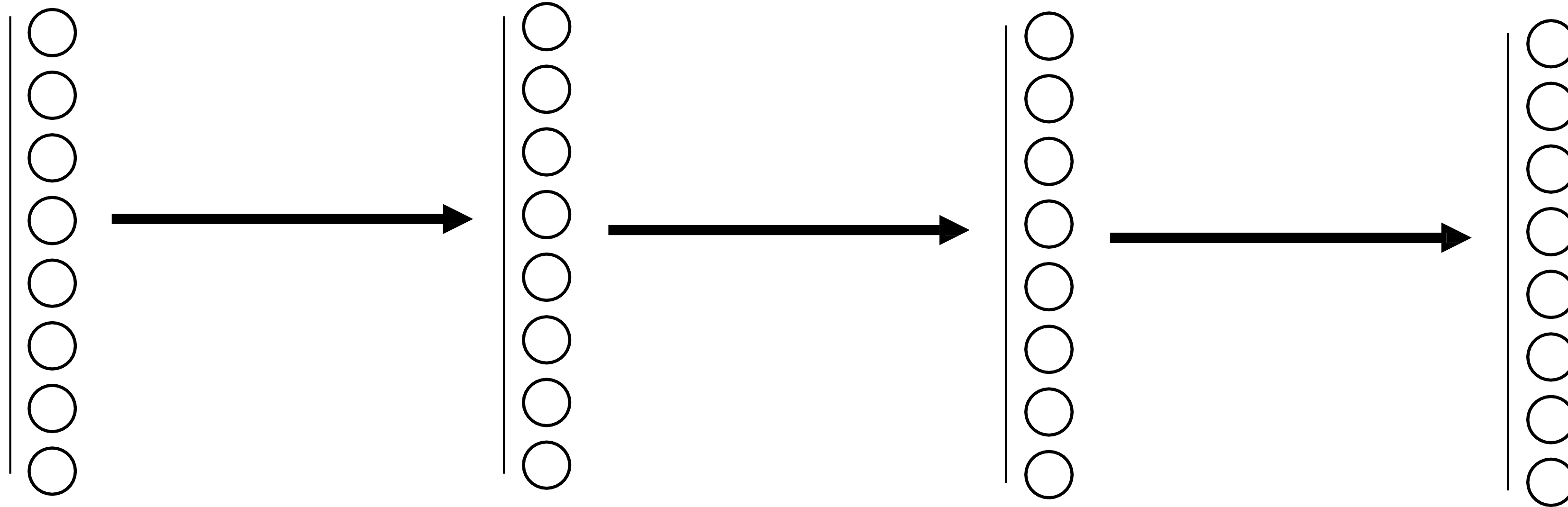
$$\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$$

Stacking layers - What's actually happening?

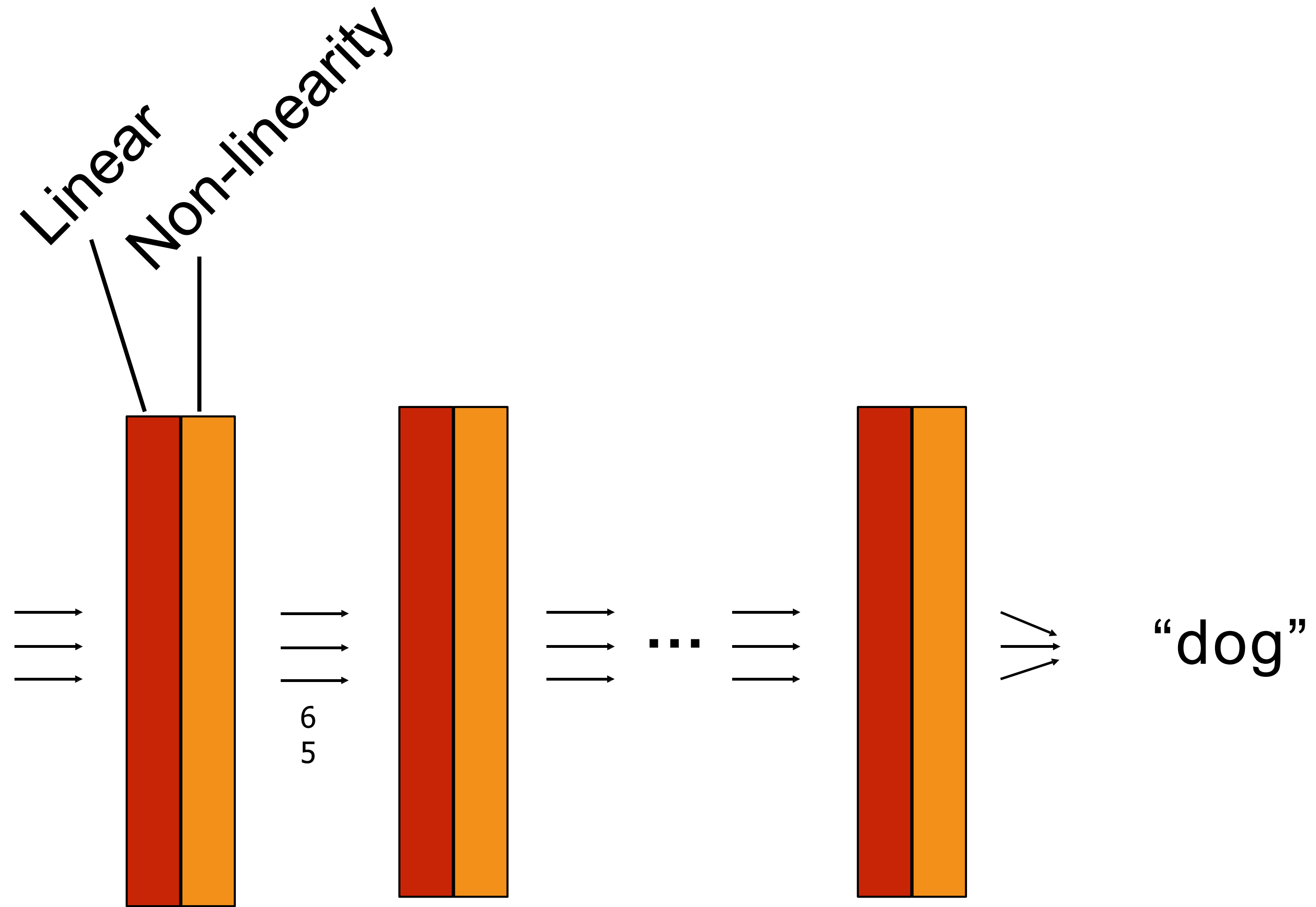
Low level features:
e.g., edge, texture, ...

higher level features:
e.g., shape

even higher level features:
e.g., "paw", "fur"



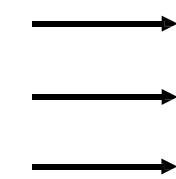
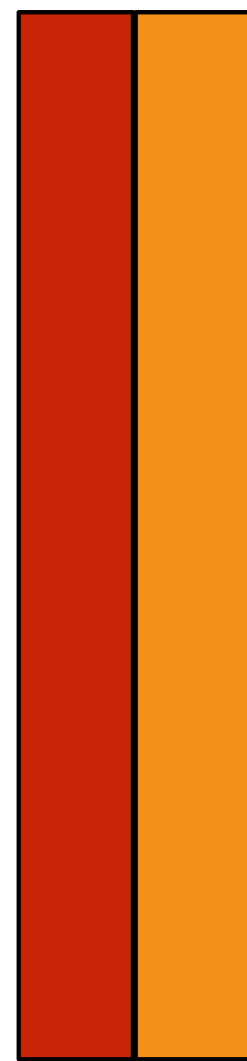
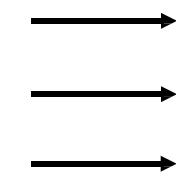
Deep nets



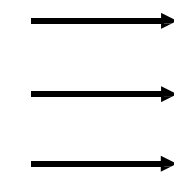
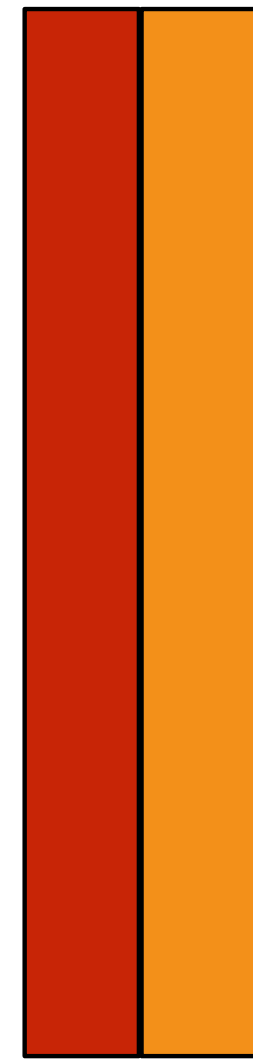
$$f(x) = f_L(\dots f_3(f_2(f_1(x))))$$

Deep nets - Intuition

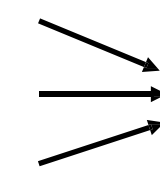
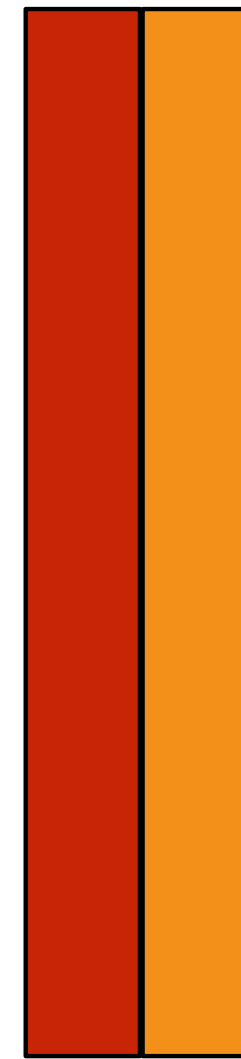
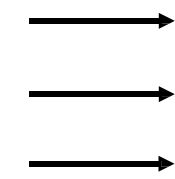
“has horizontal edge”
“has vertical edge”



6
6



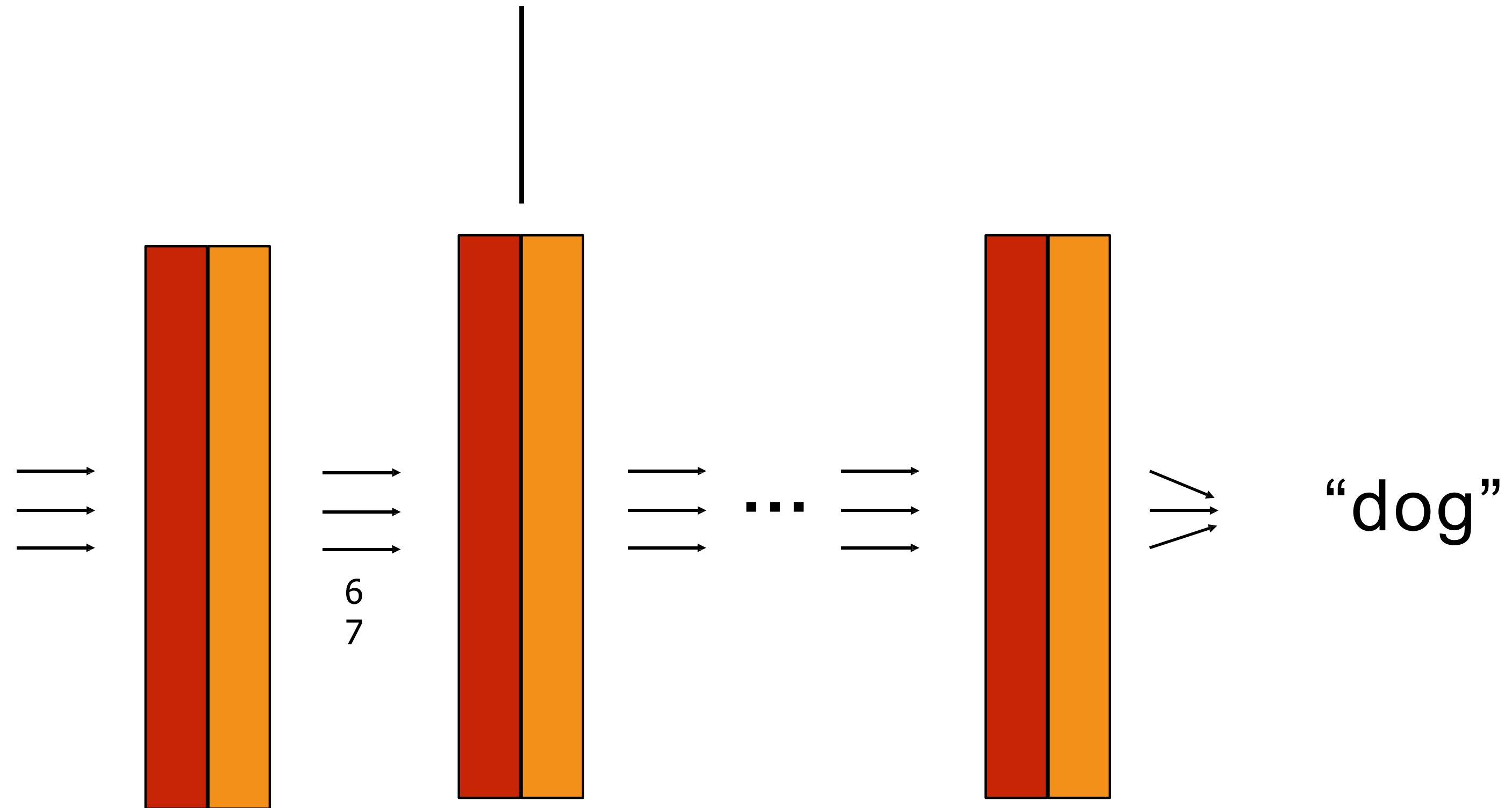
...



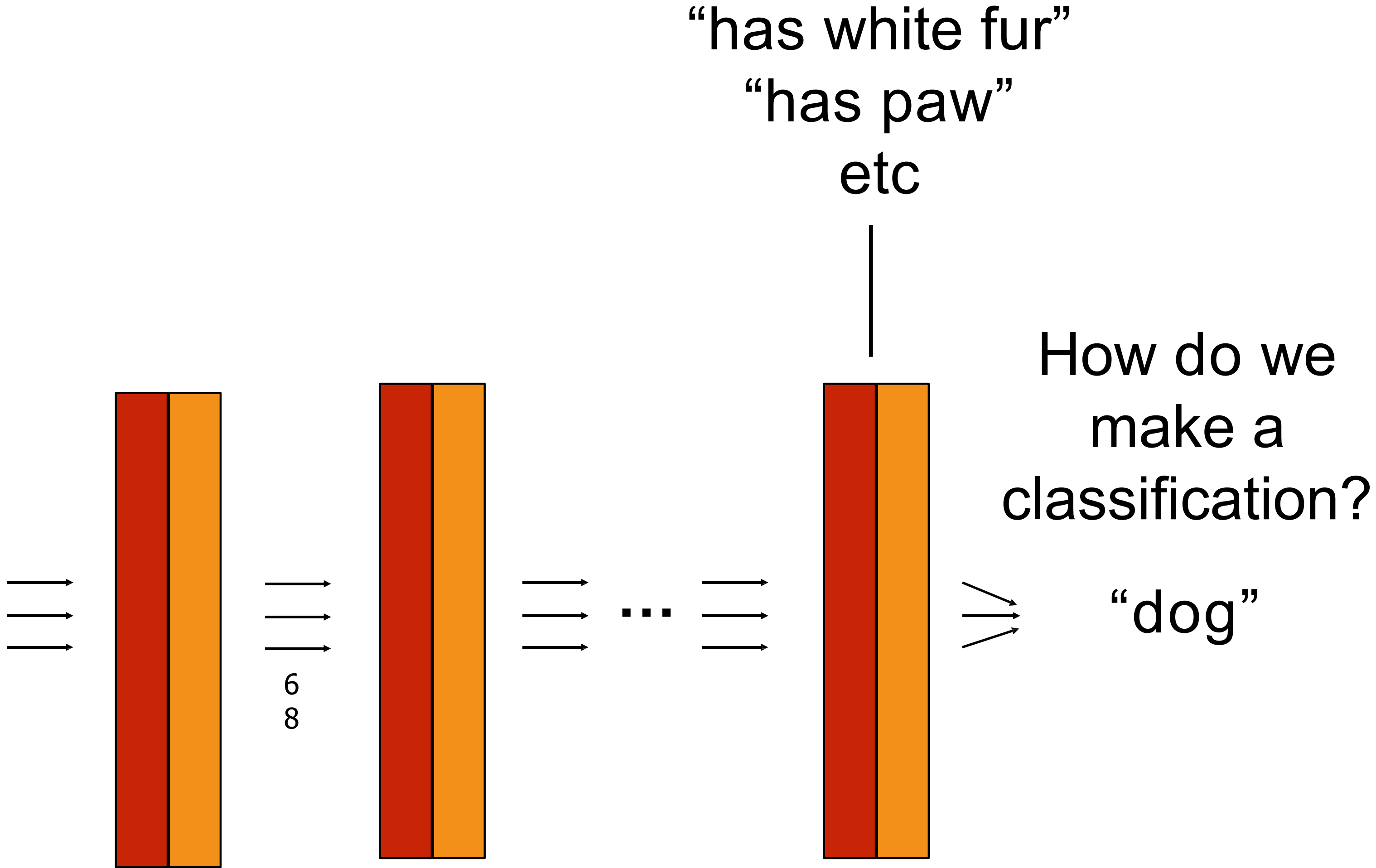
“dog”

Deep nets - Intuition

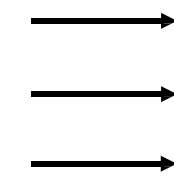
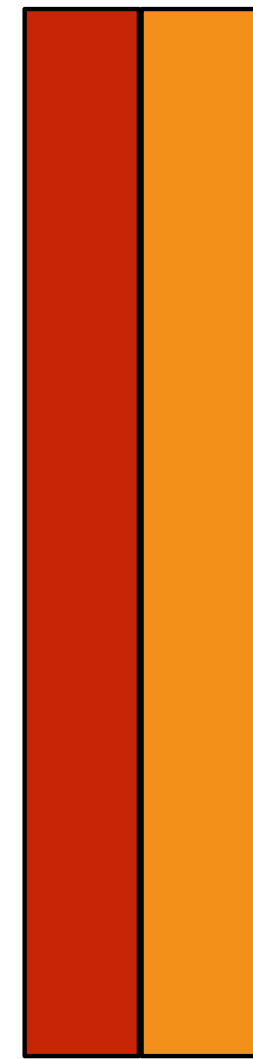
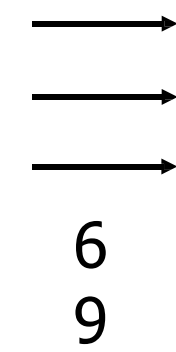
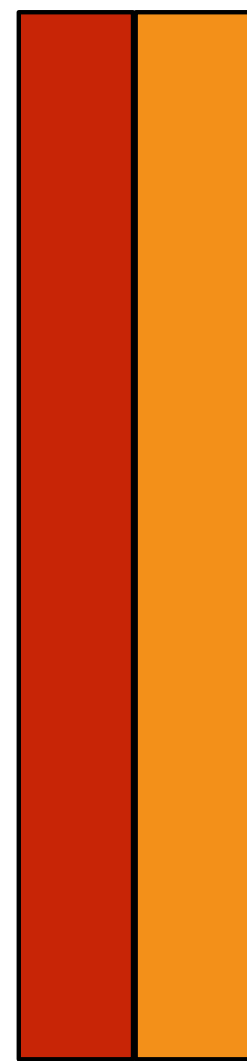
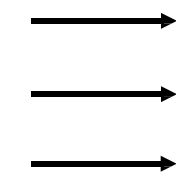
“has rounded edge”



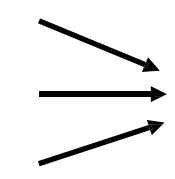
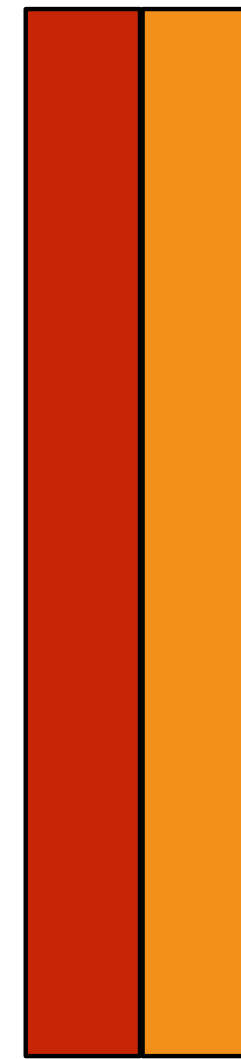
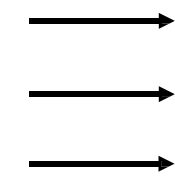
Deep nets - Intuition



Deep nets - Intuition



...



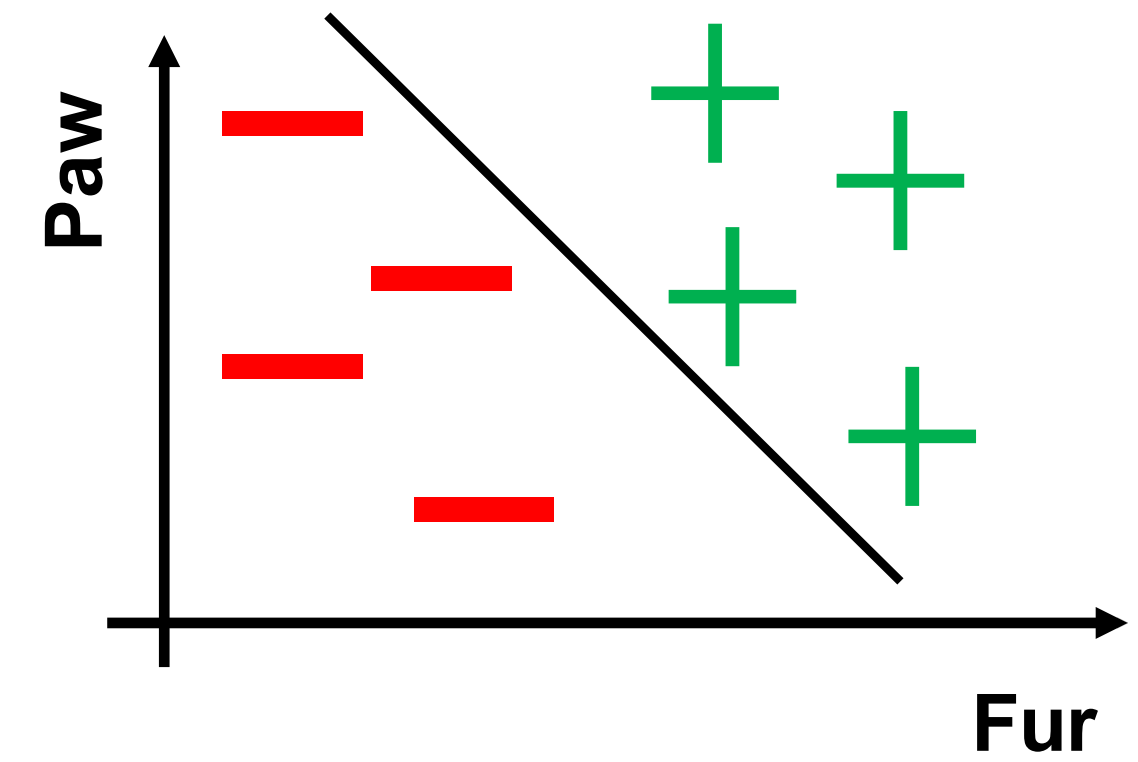
“dog”

Classify

“has white fur”
“has paw”
etc

Recall:

Feature Space



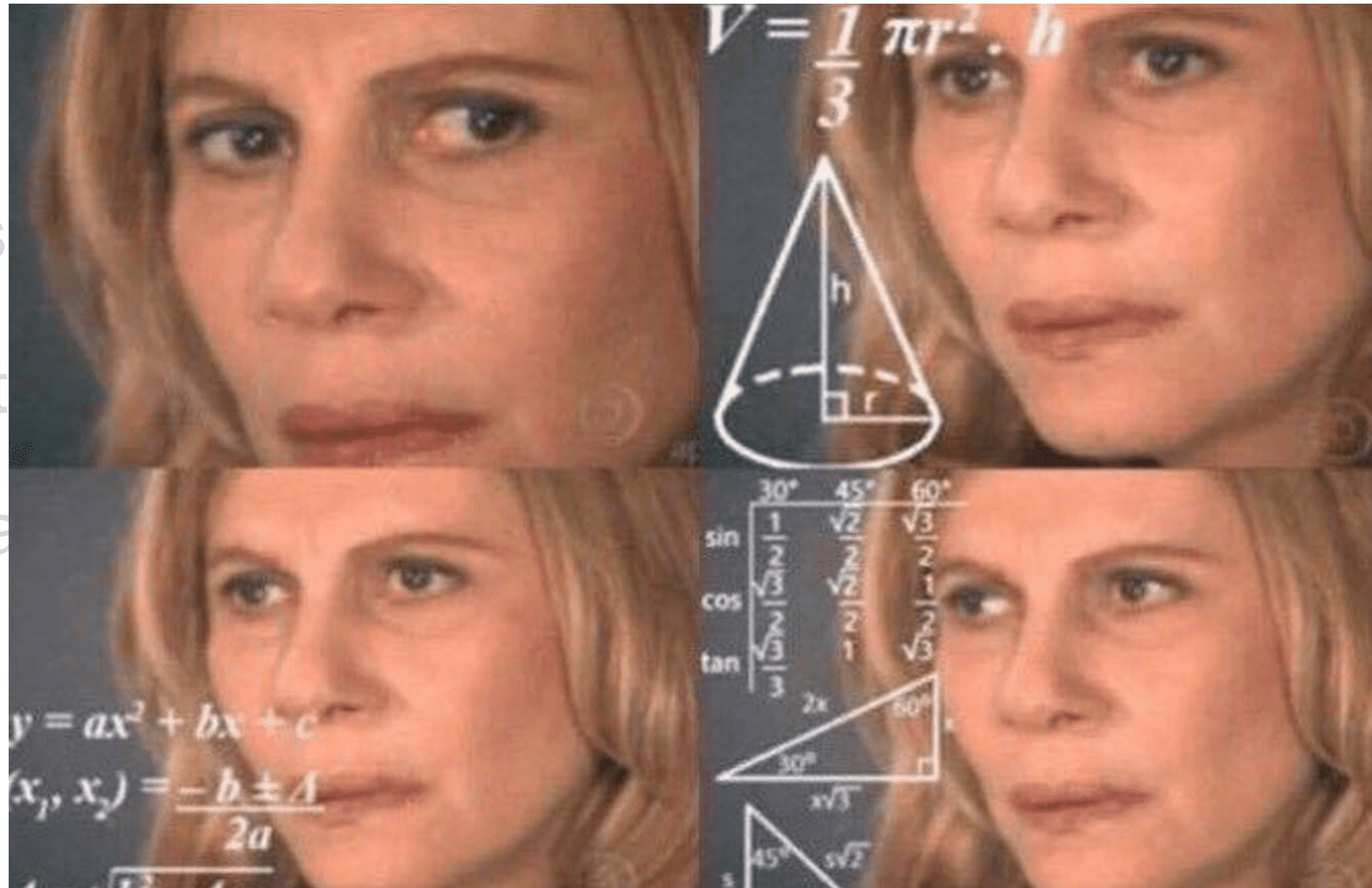
Computation has a simple form

- Composition of linear functions with nonlinearities in between
- E.g. matrix multiplications with ReLU, $\max(0, \mathbf{x})$ afterwards
- Do a matrix multiplication, set all negative values to 0, repeat

But where do we get the weights from?

Computation has a simple form

- Compos
- E.g. mat
- Do a ma



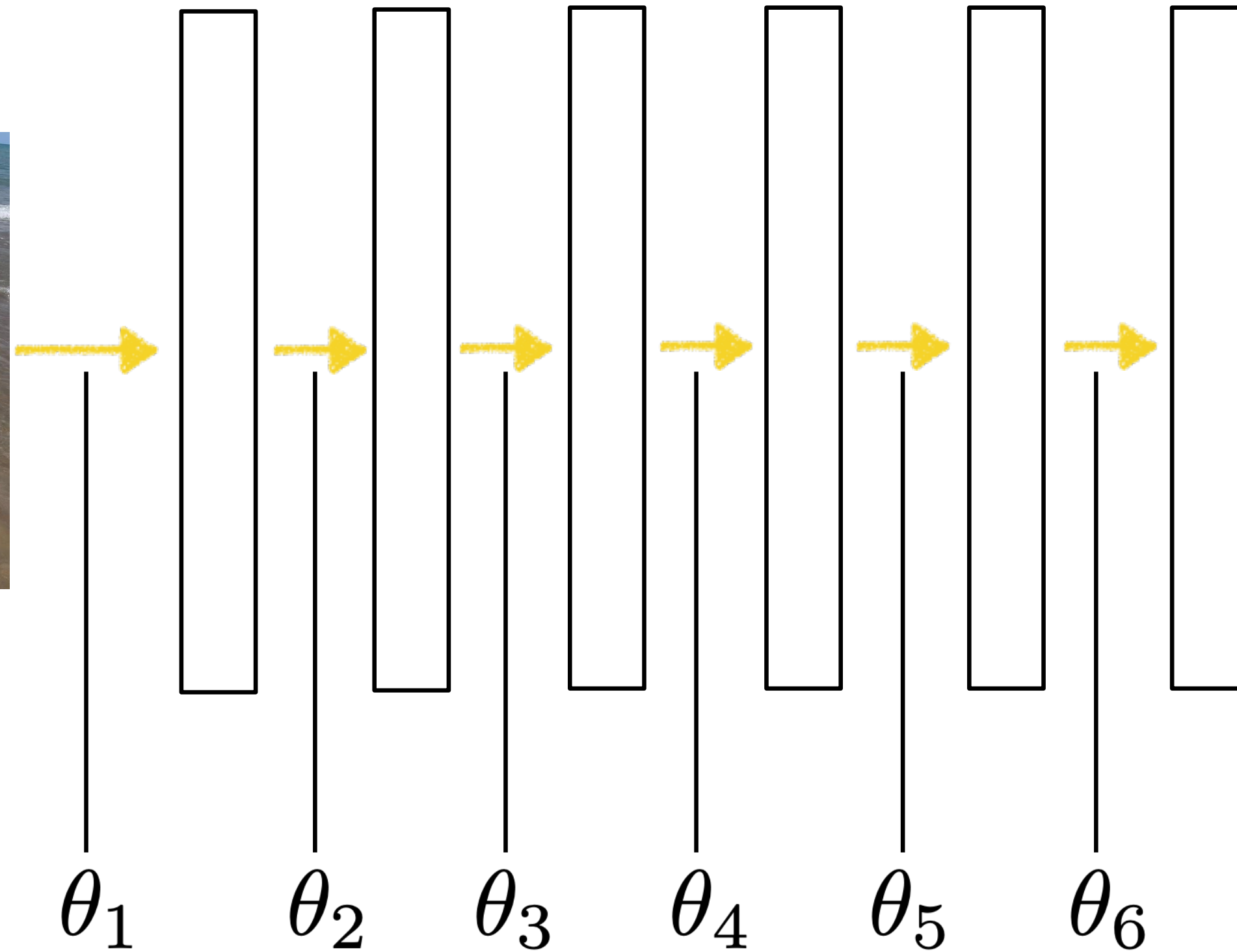
in between
afterwards
o 0, repeat

But where do we get the weights from?

How would we learn the parameters?

y_1
“dog”

x_1



$\mathcal{L}(f_{\theta}(\mathbf{x}_1), \mathbf{y}_1)$

predicted

ground truth

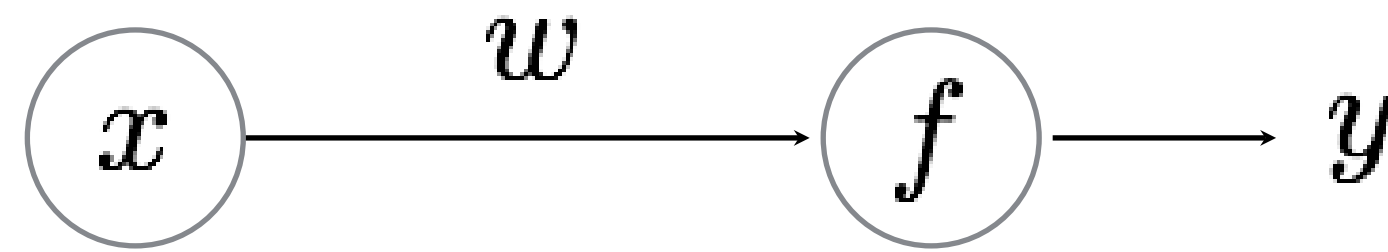
Learned

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

Training neural networks

Let's start easy

world's smallest neural network!
(a “perceptron”)



$$y = wx$$

(a.k.a. line equation, linear regression)

Training a Neural Network

Given a set of samples and a Perceptron

$$\{x_i, y_i\}$$

$$y = f_{\text{PER}}(x; w)$$

Estimate the parameter of the Perceptron

$$w$$

Given training data:

x	y
10	10.1
2	1.9
3.5	3.4
1	1.1

What do you think the weight parameter is?

$$y = wx$$

Given training data:

x	y
10	10.1
2	1.9
3.5	3.4
1	1.1

What do you think the weight parameter is?

$$y = wx$$

not so obvious as the network gets more complicated so we use ...

An Incremental Learning Strategy

(gradient descent)

Given several examples

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

and a perceptron

$$\hat{y} = wx$$

An Incremental Learning Strategy

(gradient descent)

Given several examples

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

and a perceptron

$$\hat{y} = wx$$

Modify weight w such that \hat{y} gets **'closer'** to y

An Incremental Learning Strategy

(gradient descent)

Given several examples

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

and a perceptron

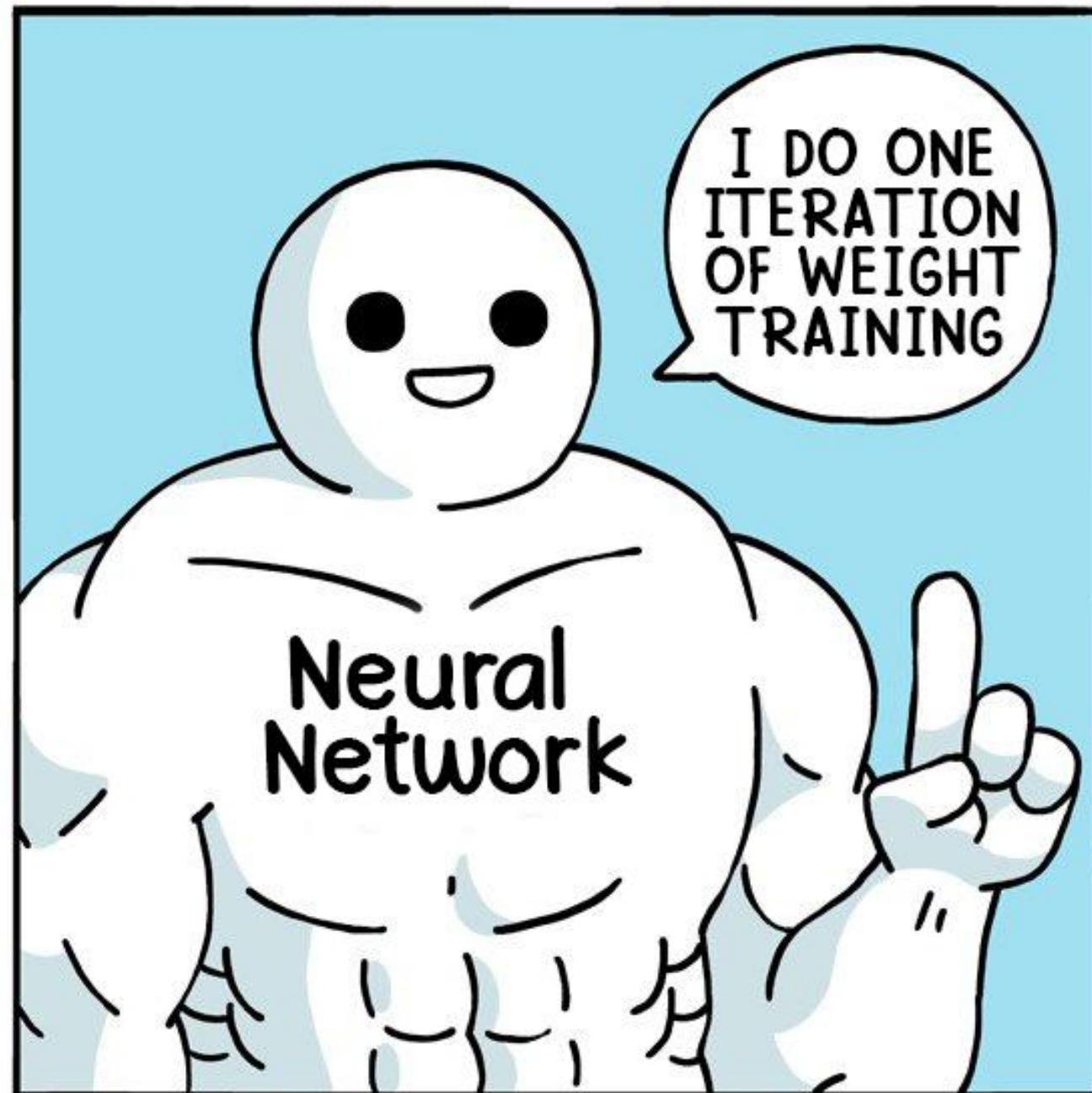
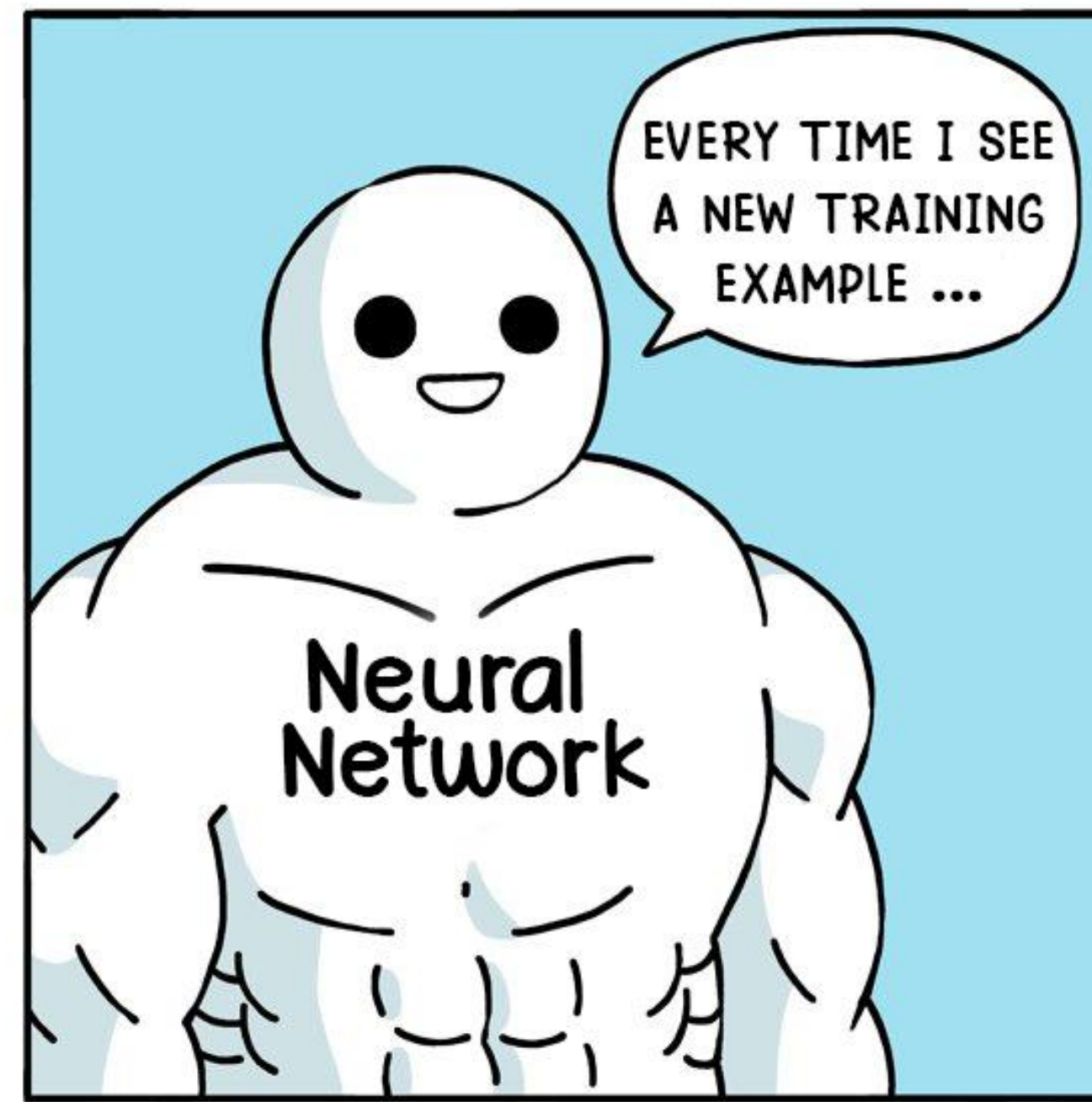
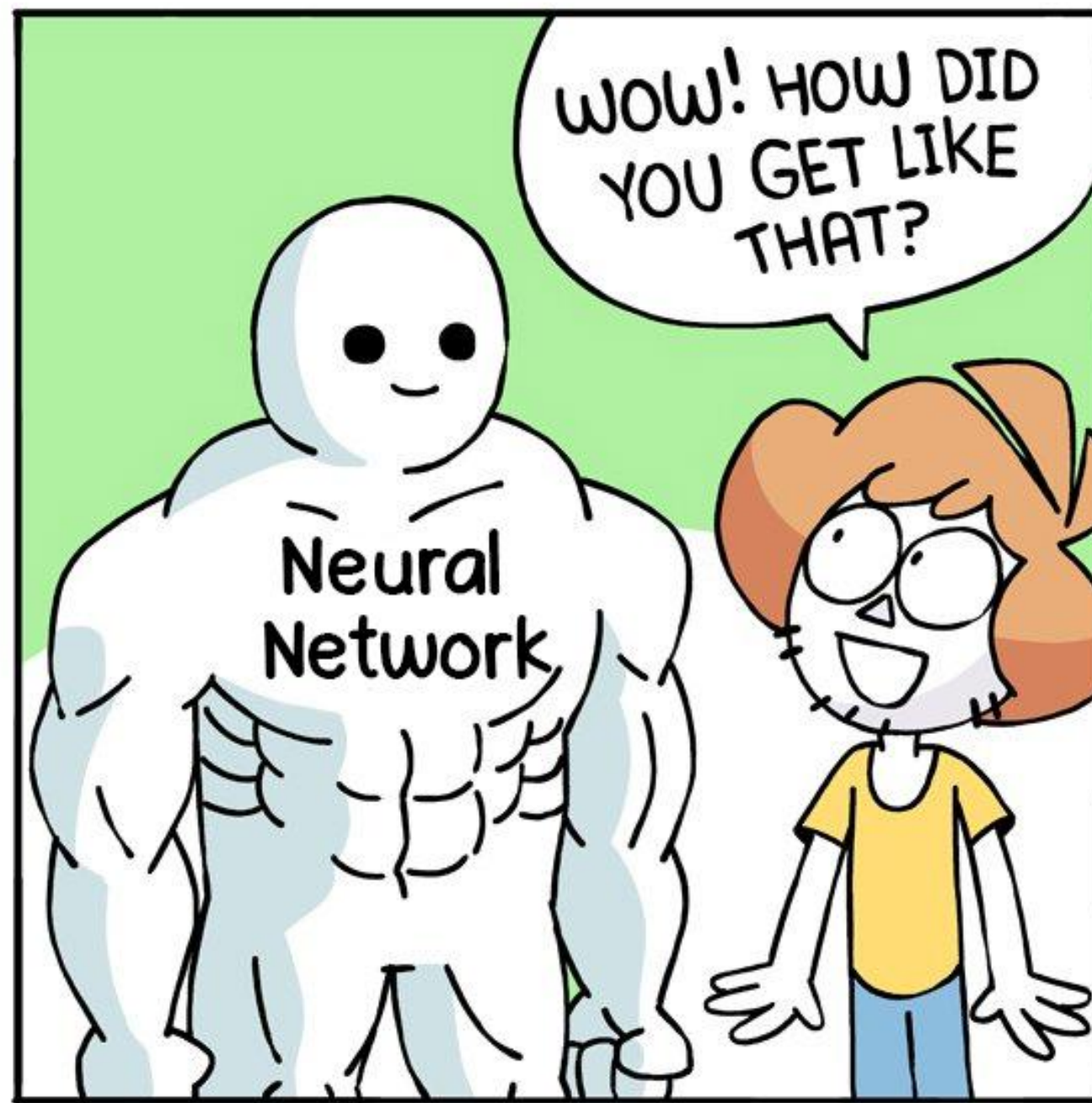
$$\hat{y} = wx$$

Modify weight w such that \hat{y} gets '**closer**' to y

perceptron
parameter

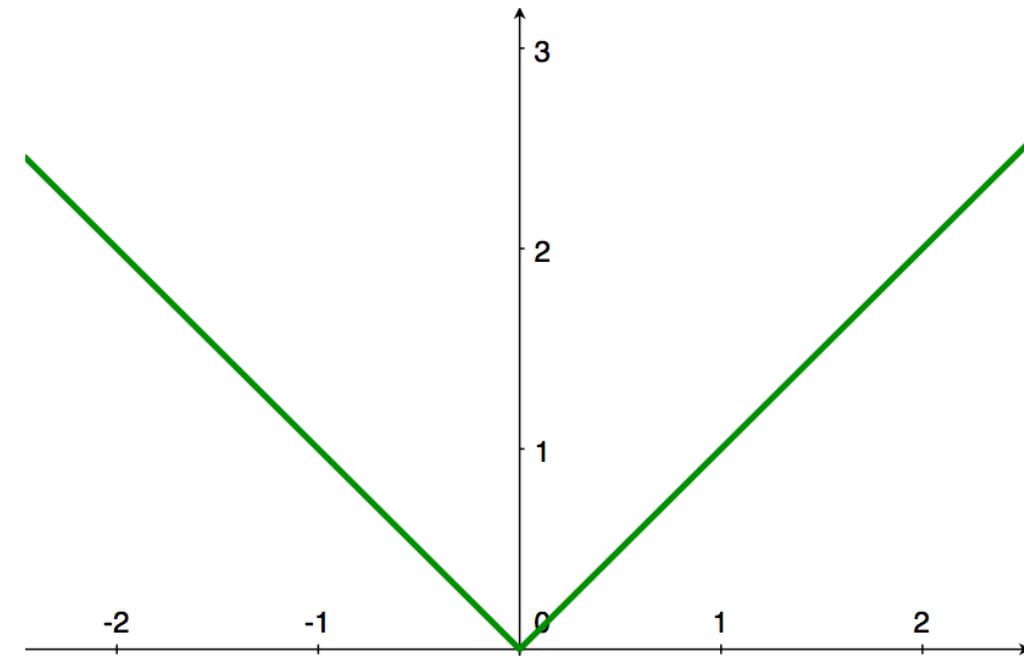
perceptron
output

true
label



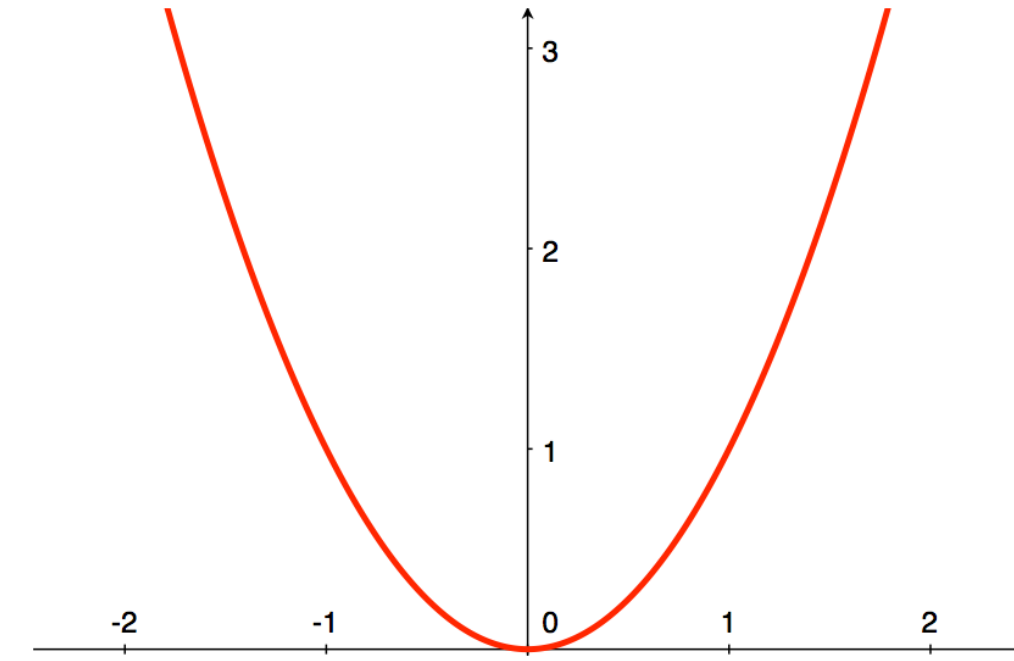
L1 Loss

$$\ell(\hat{y}, y) = |\hat{y} - y|$$



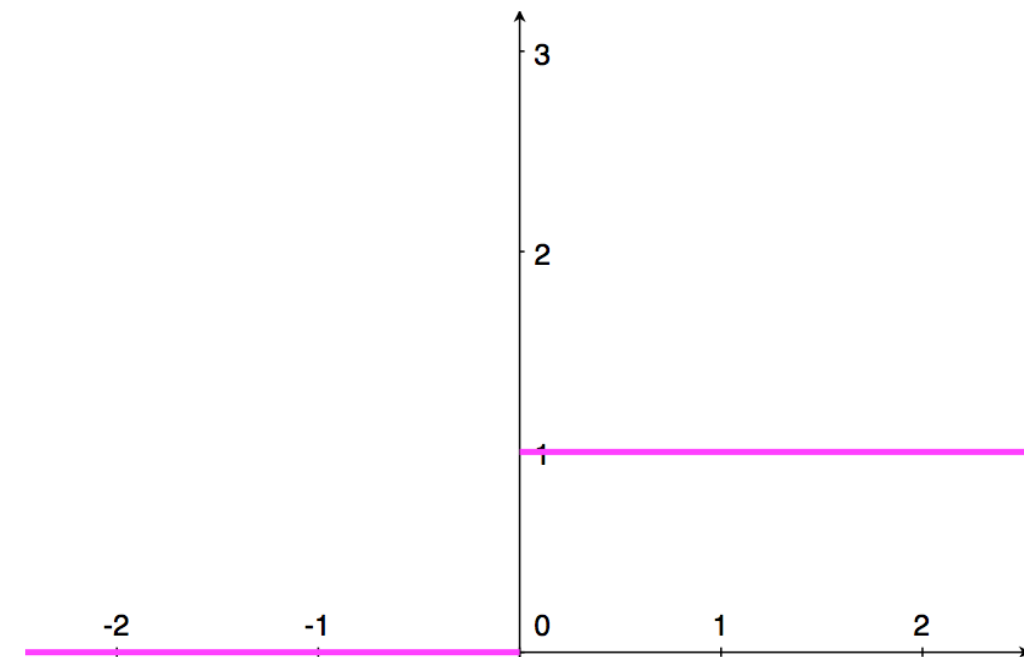
L2 Loss

$$\ell(\hat{y}, y) = (\hat{y} - y)^2$$



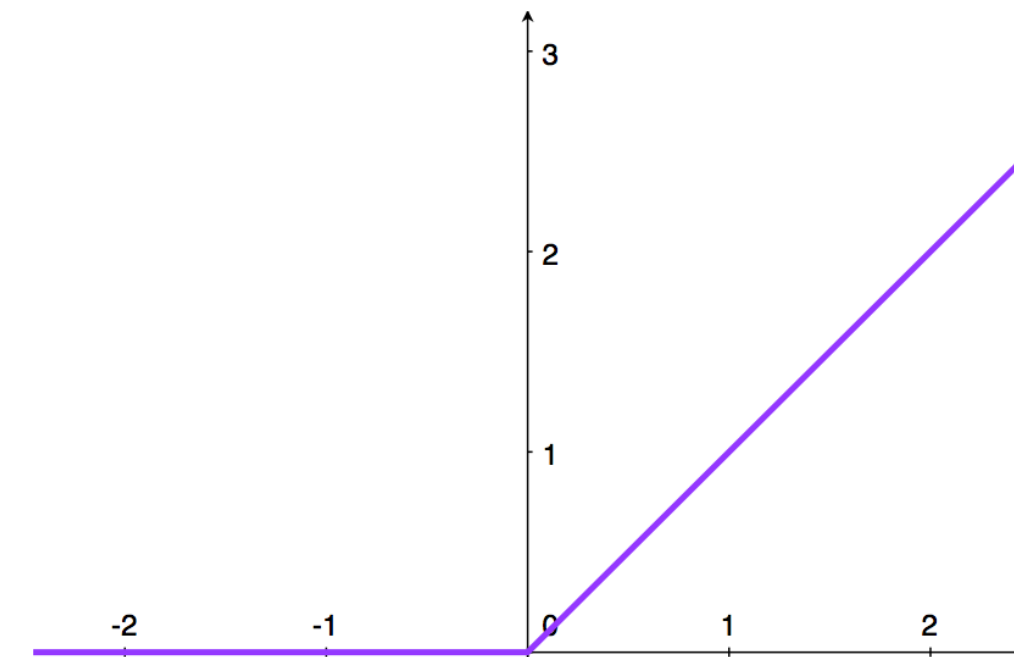
Zero-One Loss

$$\ell(\hat{y}, y) = \mathbf{1}[\hat{y} \neq y]$$

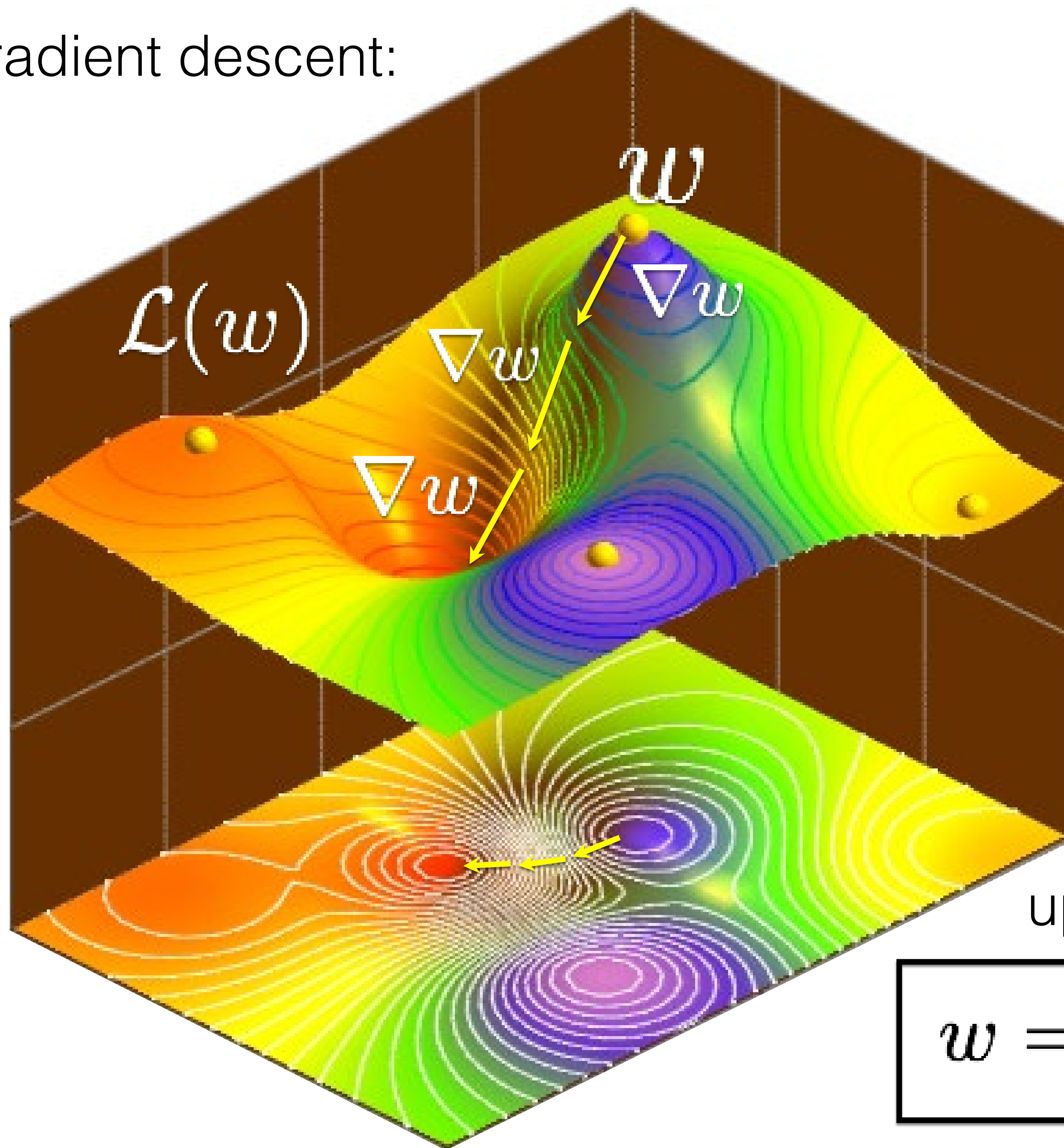


Hinge Loss

$$\ell(\hat{y}, y) = \max(0, 1 - y \cdot \hat{y})$$



Gradient descent:



update rule:

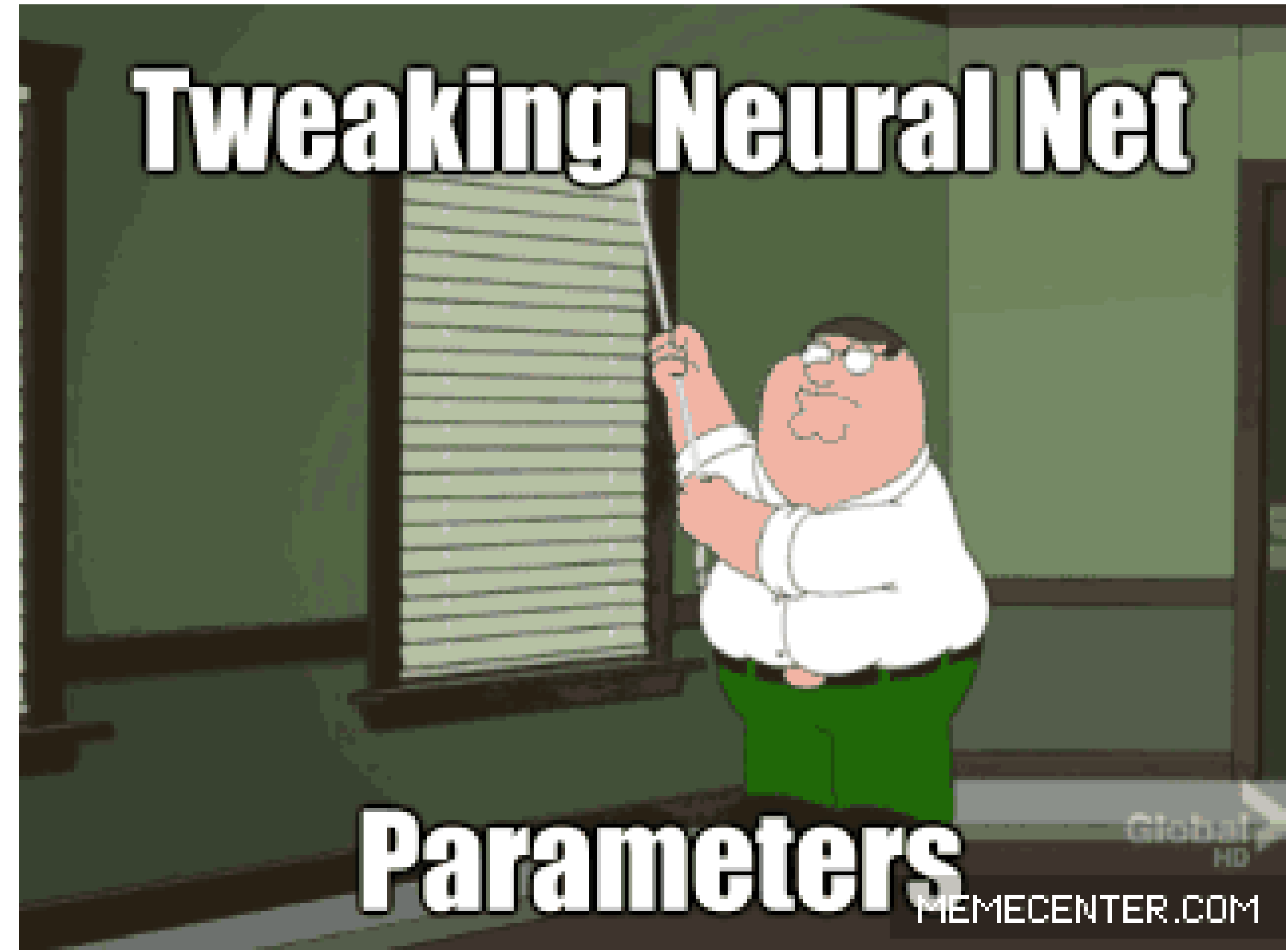
$$w = w - \nabla w$$

Backpropagation

Geoff Hinton after writing the paper on backprop in 1986



Backpropagation



$\frac{d\mathcal{L}}{dw}$...is the rate at which **this** will change...

$$\mathcal{L} = \frac{1}{2}(y - \hat{y})^2$$

the loss function

... per unit change of **this**

$$y = wx$$

the weight parameter

Let's compute the derivative...

Compute the derivative

$$\begin{aligned}\frac{d\mathcal{L}}{dw} &= \frac{d}{dw} \left\{ \frac{1}{2} (y - \hat{y})^2 \right\} \\ &= -(y - \hat{y}) \frac{dw x}{dw} \\ &= -(y - \hat{y}) x = \nabla w\end{aligned}$$

That means the weight update for **gradient descent** is:

$$\begin{aligned}w &= w - \nabla w \quad \text{move in direction of negative gradient} \\ &= w + (y - \hat{y}) x\end{aligned}$$

Gradient Descent (world's smallest perceptron)

For each sample

$$\{x_i, y_i\}$$

1. Predict

a. Forward pass

$$\hat{y} = wx_i$$

b. Compute Loss

$$\mathcal{L}_i = \frac{1}{2}(y_i - \hat{y})^2$$

2. Update

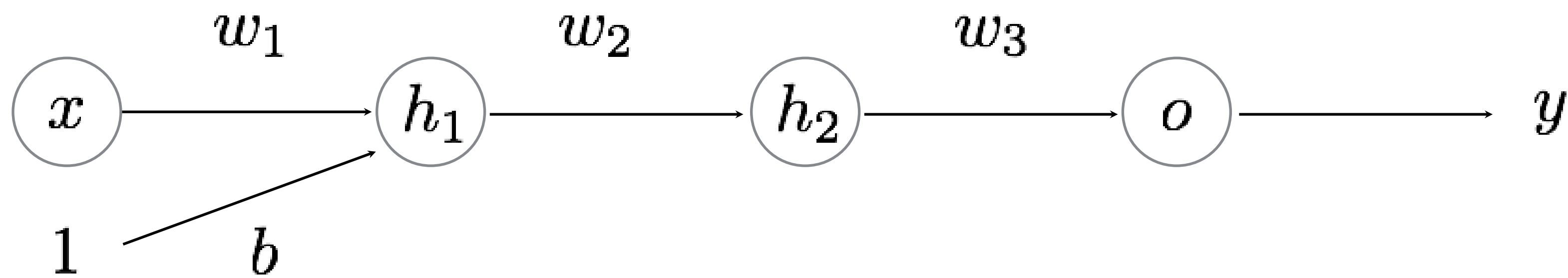
a. Back Propagation

$$\frac{d\mathcal{L}_i}{dw} = -(y_i - \hat{y})x_i = \nabla w$$

b. Gradient update

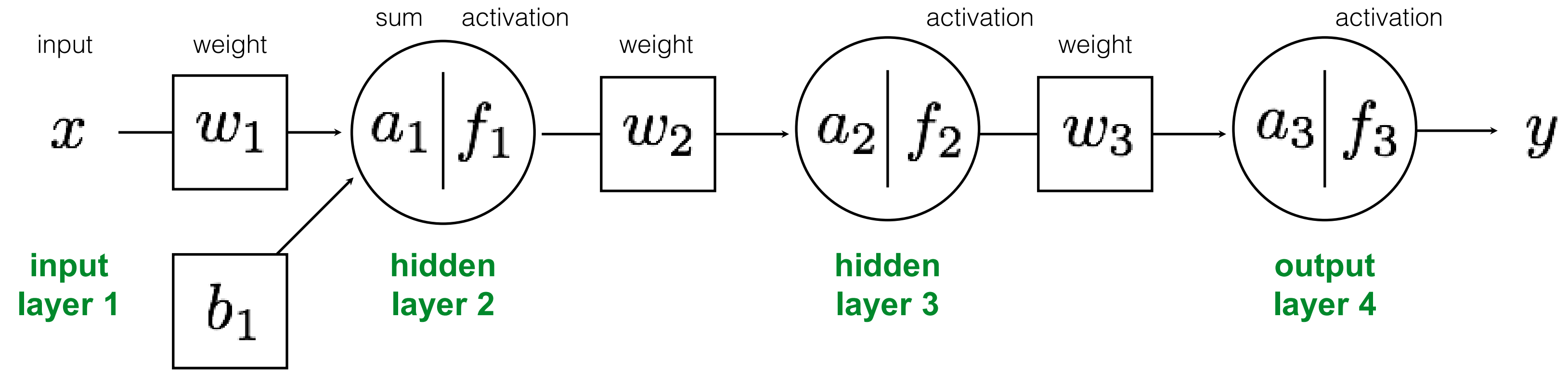
$$w = w - \nabla w$$

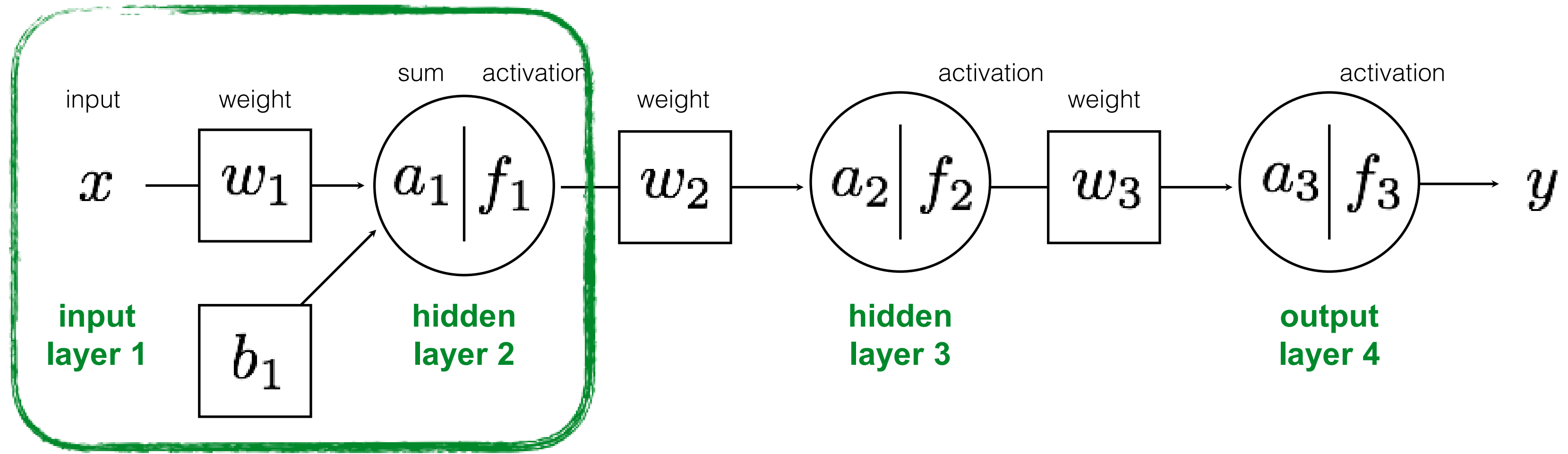
multi-layer perceptron

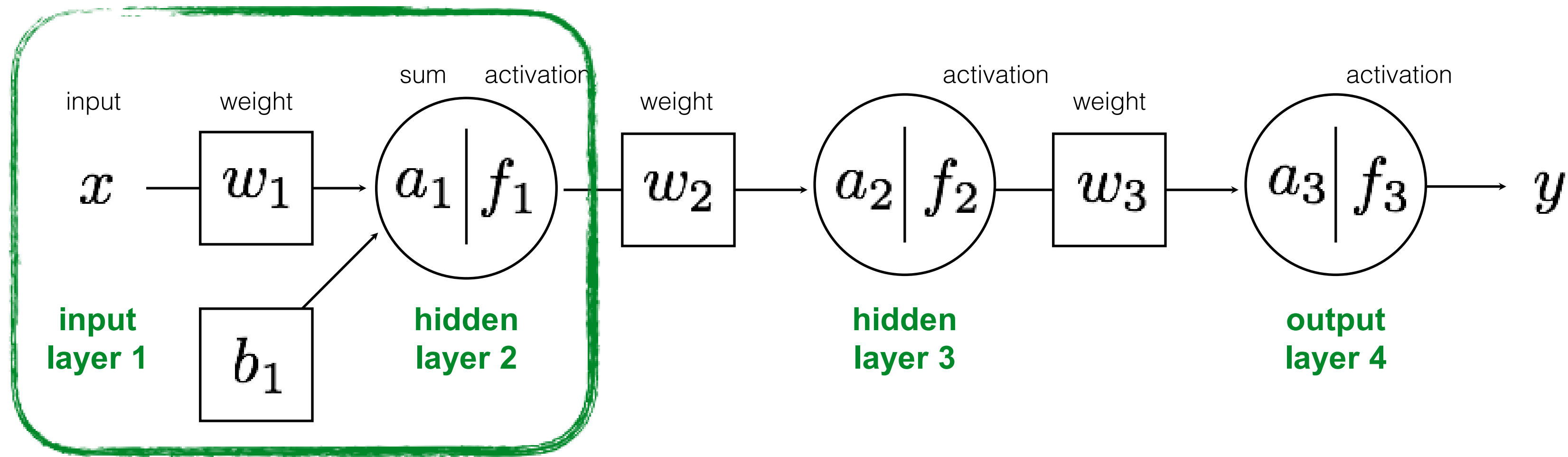


function of **FOUR** parameters and **FOUR** layers!

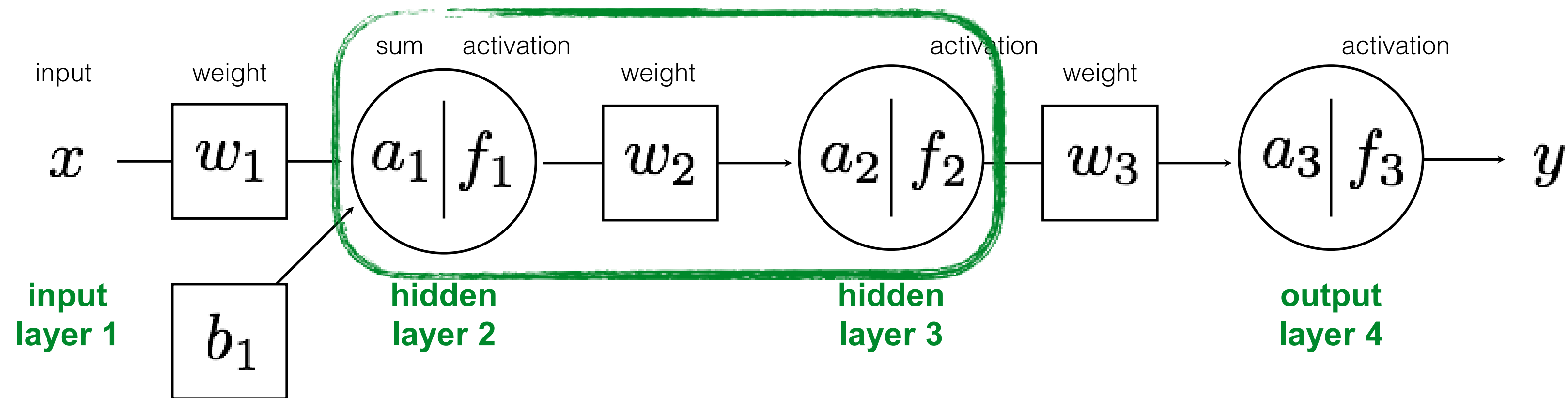




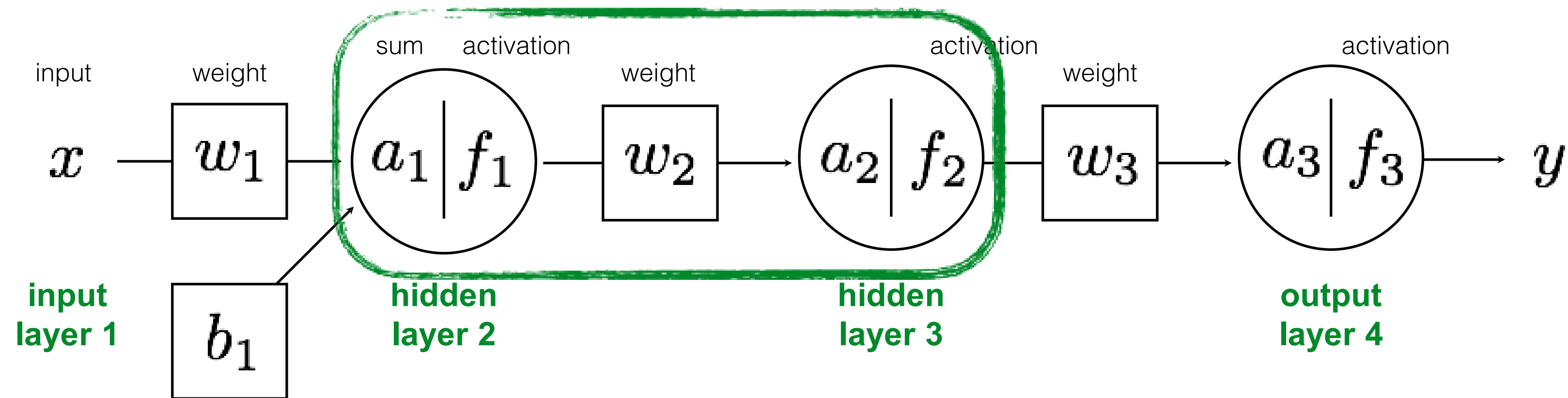




$$a_1 = w_1 \cdot x + b_1$$

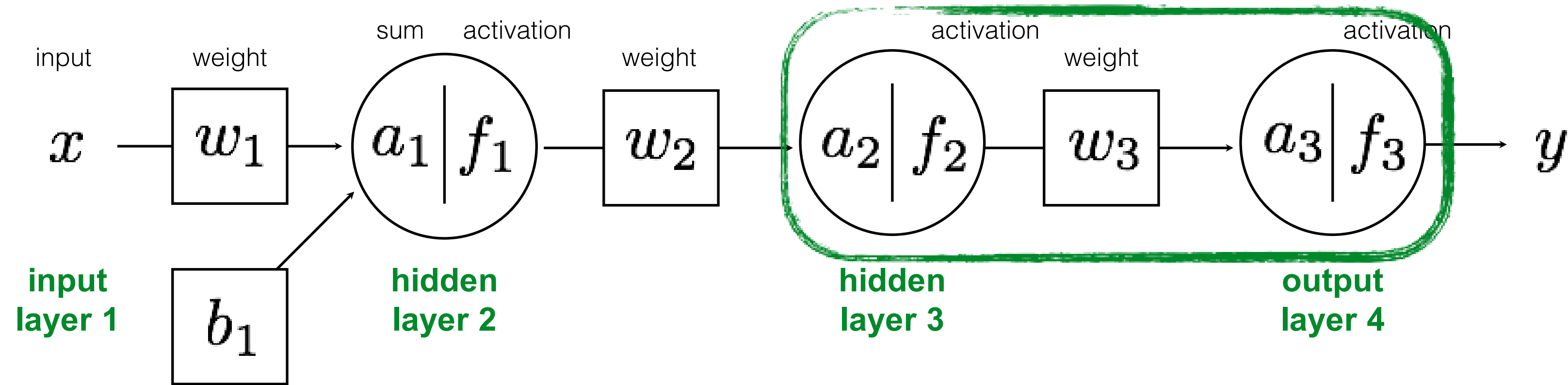


$$a_1 = w_1 \cdot x + b_1$$



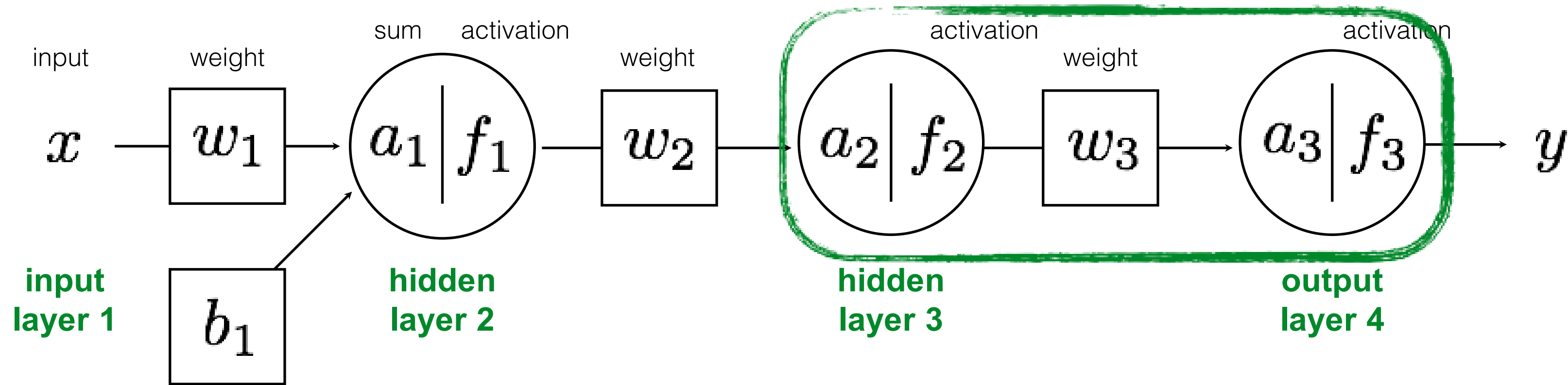
$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$



$$a_1 = w_1 \cdot x + b_1$$

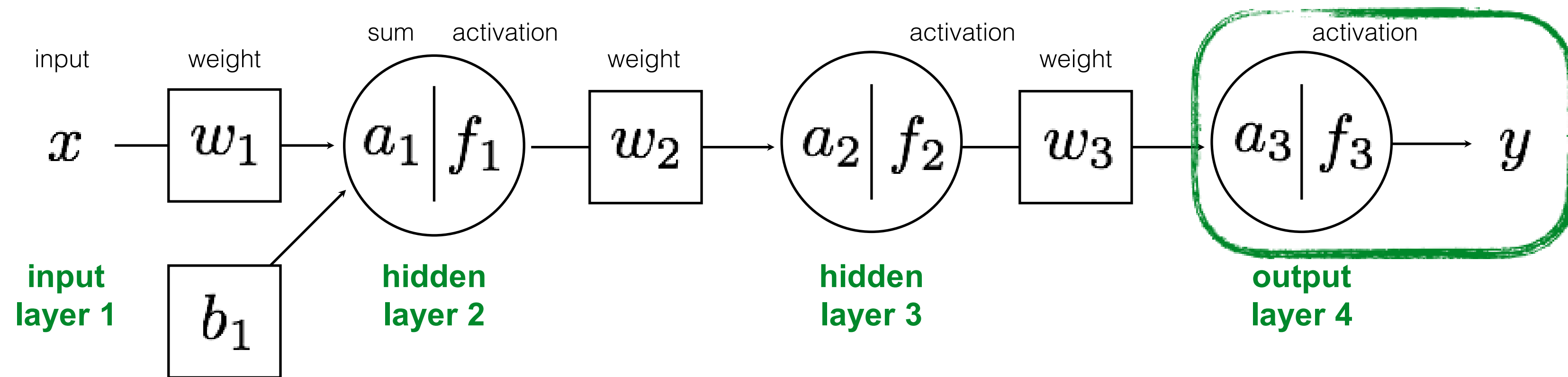
$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$



$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$

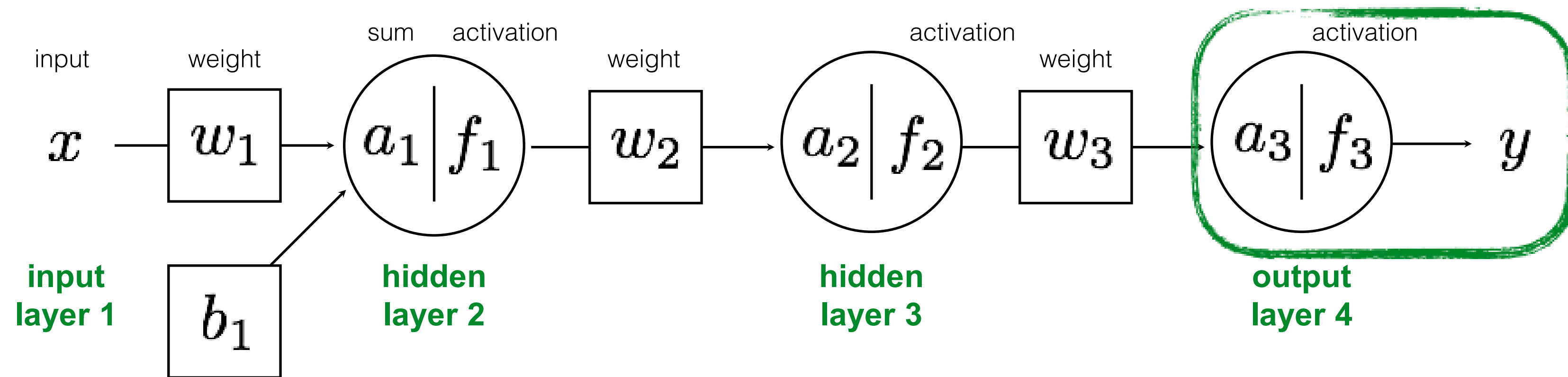
$$a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$$



$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$

$$a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$$



$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$

$$a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$$

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

Entire network can be written out as one long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

We need to train the network:

What is known? What is unknown?

Entire network can be written out as a long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$



known

We need to train the network:

What is known? What is unknown?

Entire network can be written out as a long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

activation function
sometimes has
parameters

unknown



We need to train the network:

What is known? What is unknown?

Learning an MLP

Given a set of samples and a MLP

$$\{x_i, y_i\}$$

$$y = f_{\text{MLP}}(x; \theta)$$

Estimate the parameters of the MLP

$$\theta = \{f, w, b\}$$

Gradient Descent

For each **random** sample $\{x_i, y_i\}$

1. Predict

a. Forward pass

$$\hat{y} = f_{\text{MLP}}(x_i; \theta)$$

b. Compute Loss

2. Update

a. Back Propagation

$$\frac{\partial \mathcal{L}}{\partial \theta}$$

vector of parameter partial derivatives

b. Gradient update

$$\theta \leftarrow \theta - \eta \nabla \theta$$

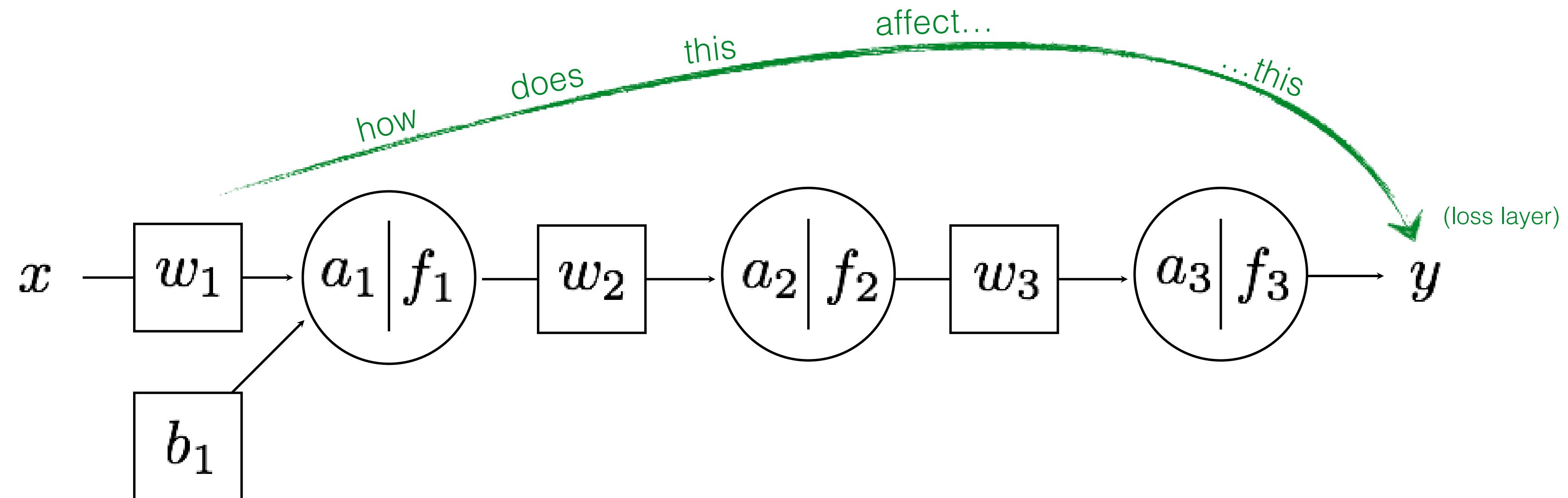
vector of parameter update equations

So we need to compute the partial derivatives

$$\frac{\partial \mathcal{L}}{\partial \theta} = \left[\frac{\partial \mathcal{L}}{\partial w_3} \quad \frac{\partial \mathcal{L}}{\partial w_2} \quad \frac{\partial \mathcal{L}}{\partial w_1} \quad \frac{\partial \mathcal{L}}{\partial b} \right]$$

Remember,

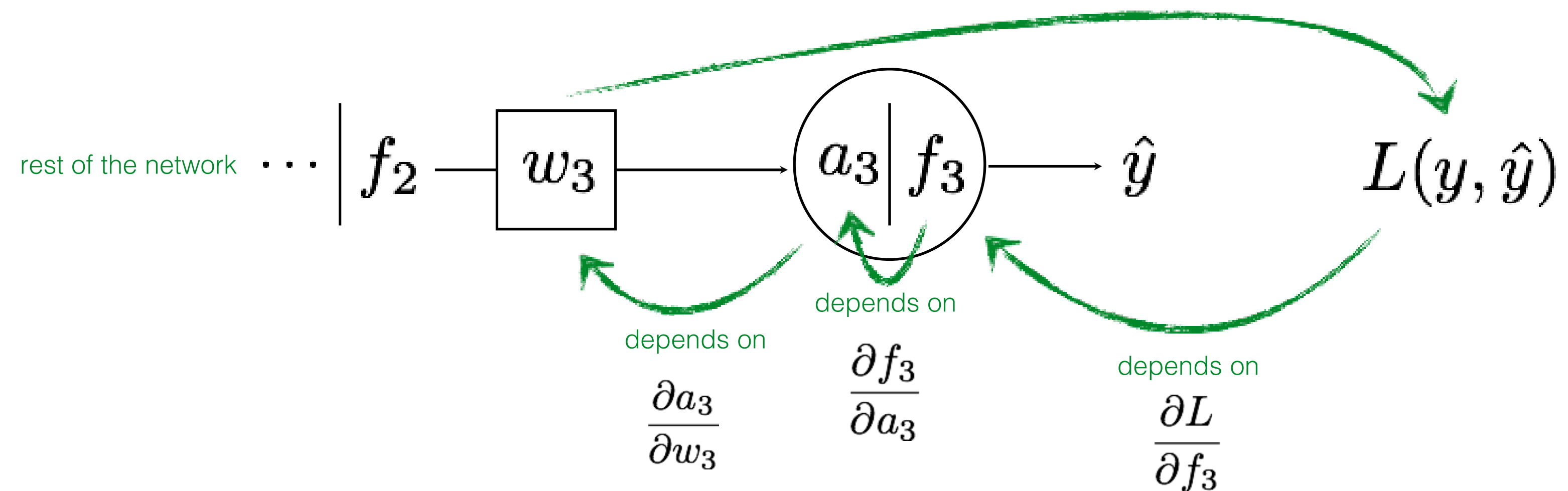
Partial derivative $\frac{\partial L}{\partial w_1}$ describes...

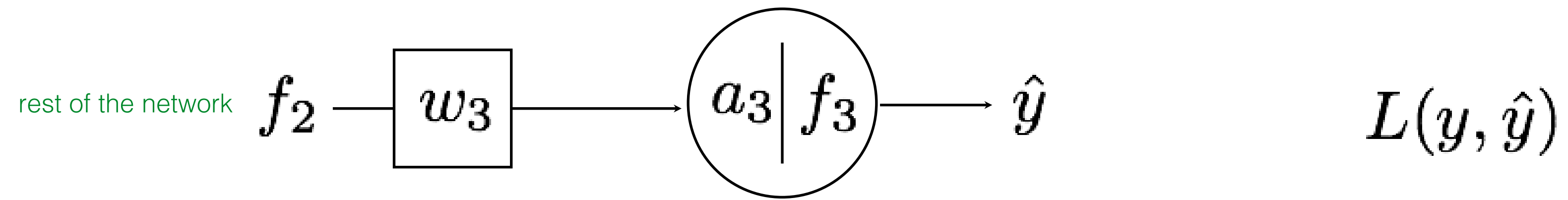


According to the chain rule...

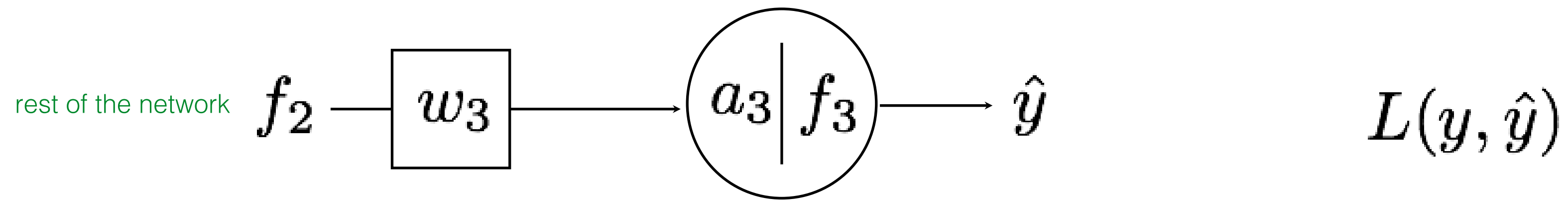
$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

Intuitively, the effect of weight on loss function : $\frac{\partial L}{\partial w_3}$





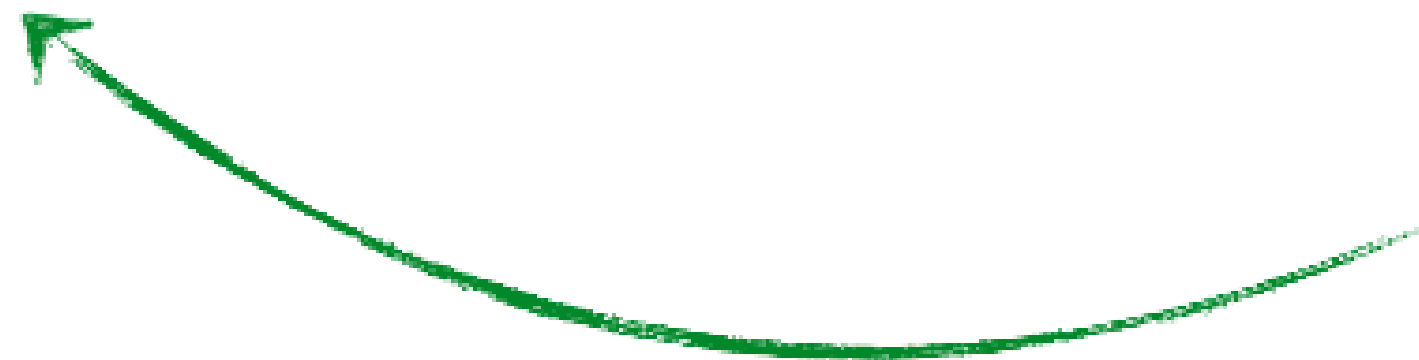
$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \quad \text{Chain Rule!}$$

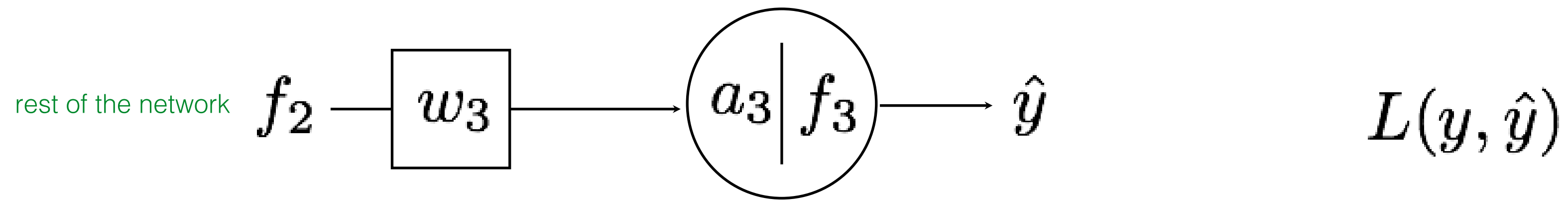


$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

$$= -\eta(y - \hat{y}) \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

Just the partial derivative of L2 loss



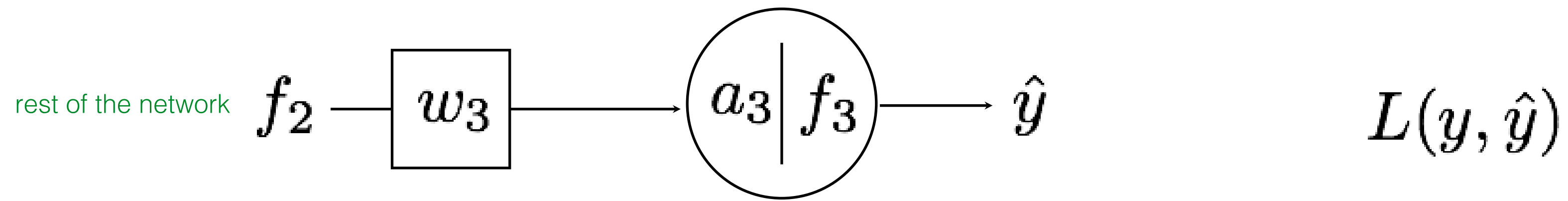


$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

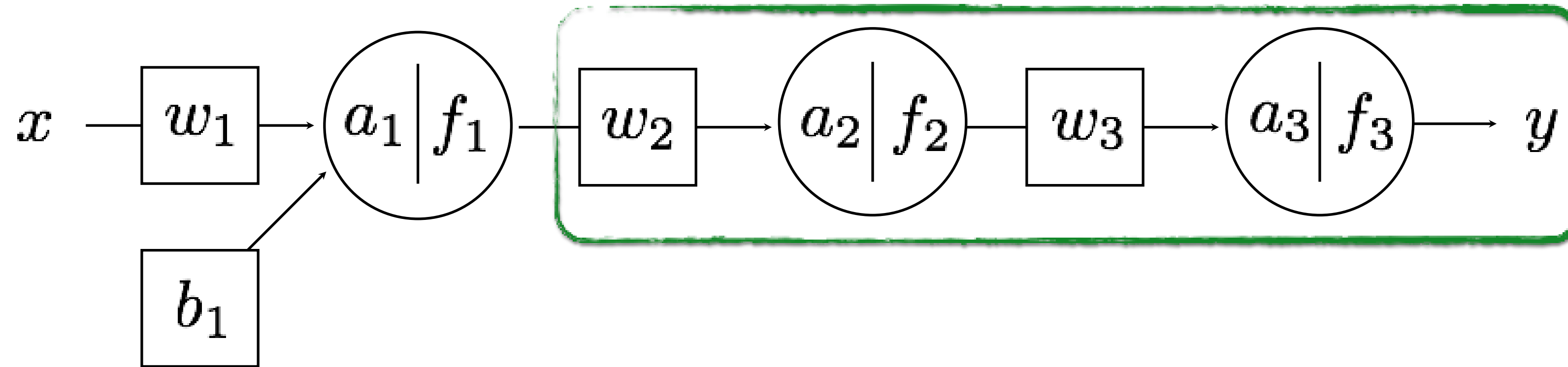
$$= -\eta(y - \hat{y}) \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

Let's use a Sigmoid function

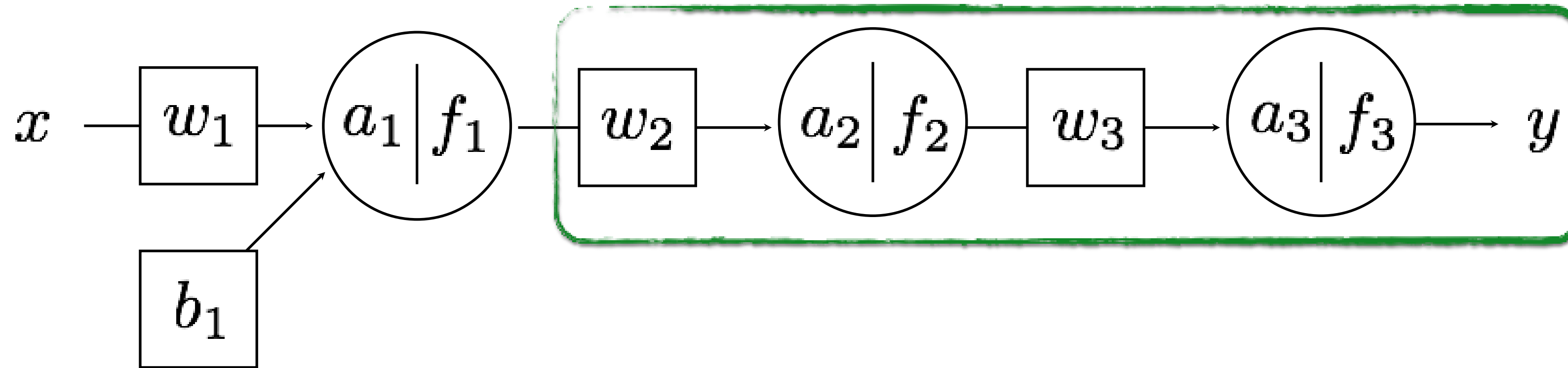
$$\frac{ds(x)}{dx} = s(x)(1 - s(x))$$



$$\begin{aligned}
 \frac{\partial L}{\partial w_3} &= \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\
 &= -\eta(y - \hat{y}) \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\
 &= -\eta(y - \hat{y}) f_3(1 - f_3) \frac{\partial a_3}{\partial w_3} \\
 &= -\eta(y - \hat{y}) f_3(1 - f_3) f_2
 \end{aligned}$$



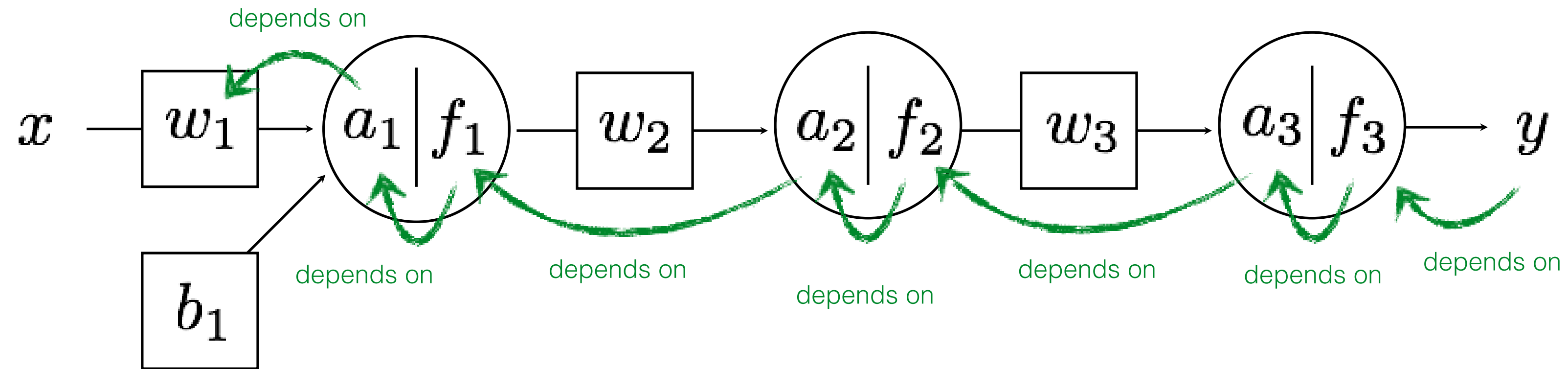
$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2}$$



$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2}$$

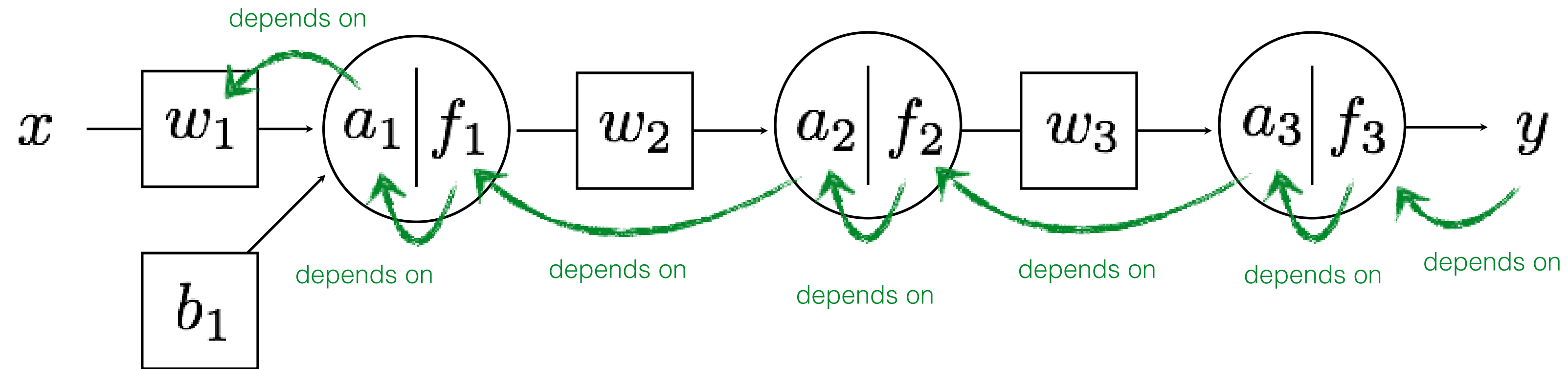
already computed.
re-use (propagate)!

The chain rule says...



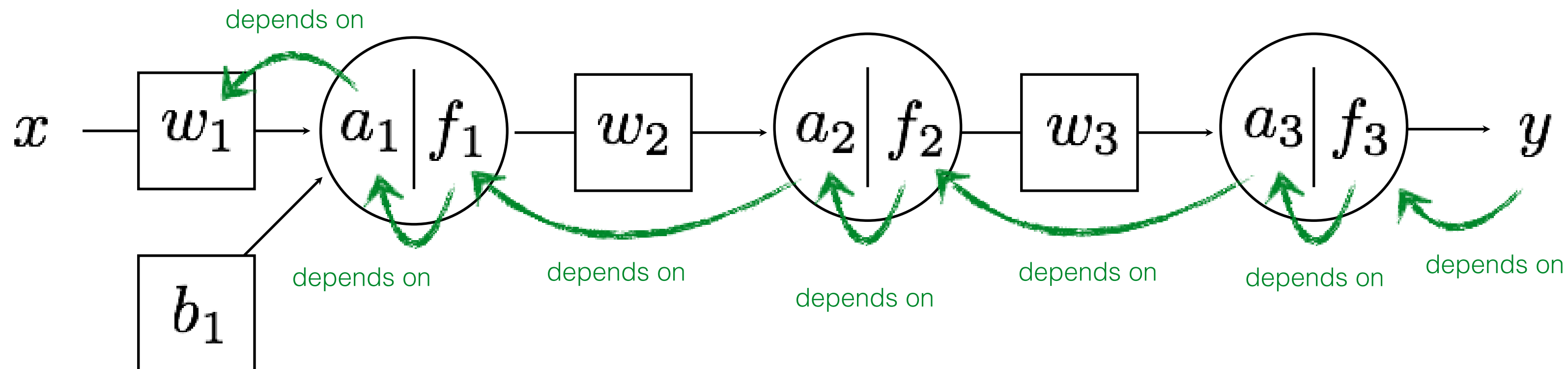
$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$

The chain rule says...

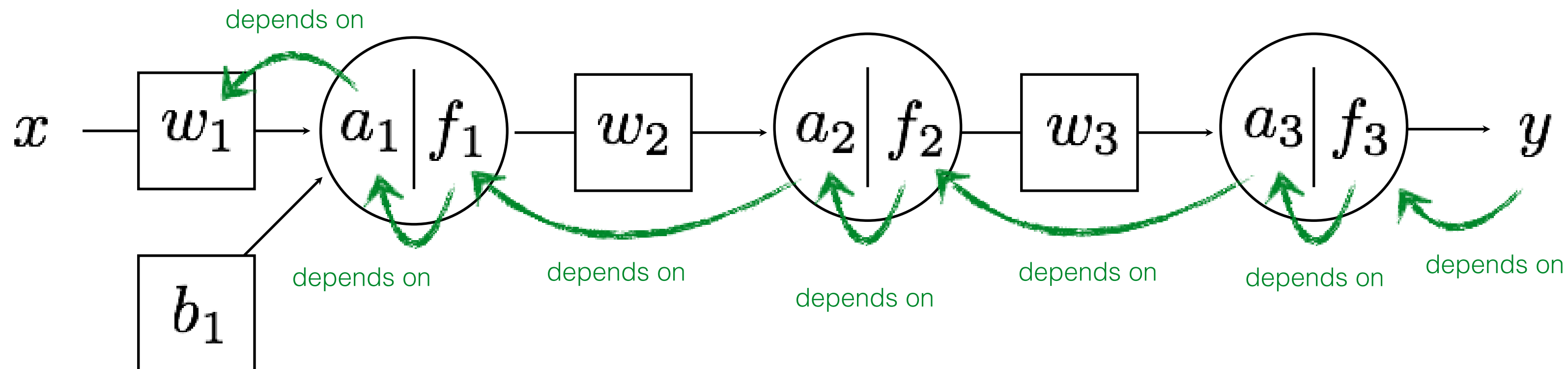


$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$

already computed.
re-use (propagate)!



$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial w_3} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\
 \frac{\partial \mathcal{L}}{\partial w_2} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2} \\
 \frac{\partial \mathcal{L}}{\partial w_1} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} \\
 \frac{\partial \mathcal{L}}{\partial b} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial b}
 \end{aligned}$$

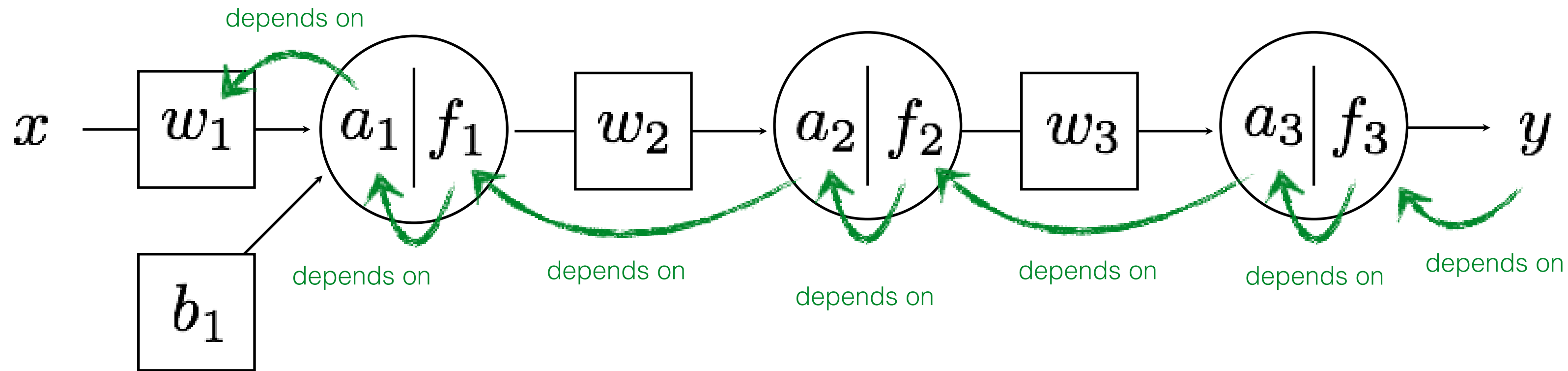


$$\frac{\partial \mathcal{L}}{\partial w_3} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2}$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial b}$$



$$\frac{\partial \mathcal{L}}{\partial w_3} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2}$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial b}$$

Gradient Descent

For each example sample $\{x_i, y_i\}$

1. Predict

a. Forward pass

$$\hat{y} = f_{\text{MLP}}(x_i; \theta)$$

b. Compute Loss

$$\mathcal{L}_i$$

2. Update

a. Back Propagation

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial w_3} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\ \frac{\partial \mathcal{L}}{\partial w_2} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2} \\ \frac{\partial \mathcal{L}}{\partial w_1} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} \\ \frac{\partial \mathcal{L}}{\partial b} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial b}\end{aligned}$$

b. Gradient update

$$w_3 = w_3 - \eta \nabla w_3$$

$$w_2 = w_2 - \eta \nabla w_2$$

$$w_1 = w_1 - \eta \nabla w_1$$

$$b = b - \eta \nabla b$$

Gradient Descent

For each example sample $\{x_i, y_i\}$

1. Predict

a. Forward pass

$$\hat{y} = f_{\text{MLP}}(x_i; \theta)$$

b. Compute Loss

$$\mathcal{L}_i$$

2. Update

a. Back Propagation

$$\frac{\partial \mathcal{L}}{\partial \theta}$$

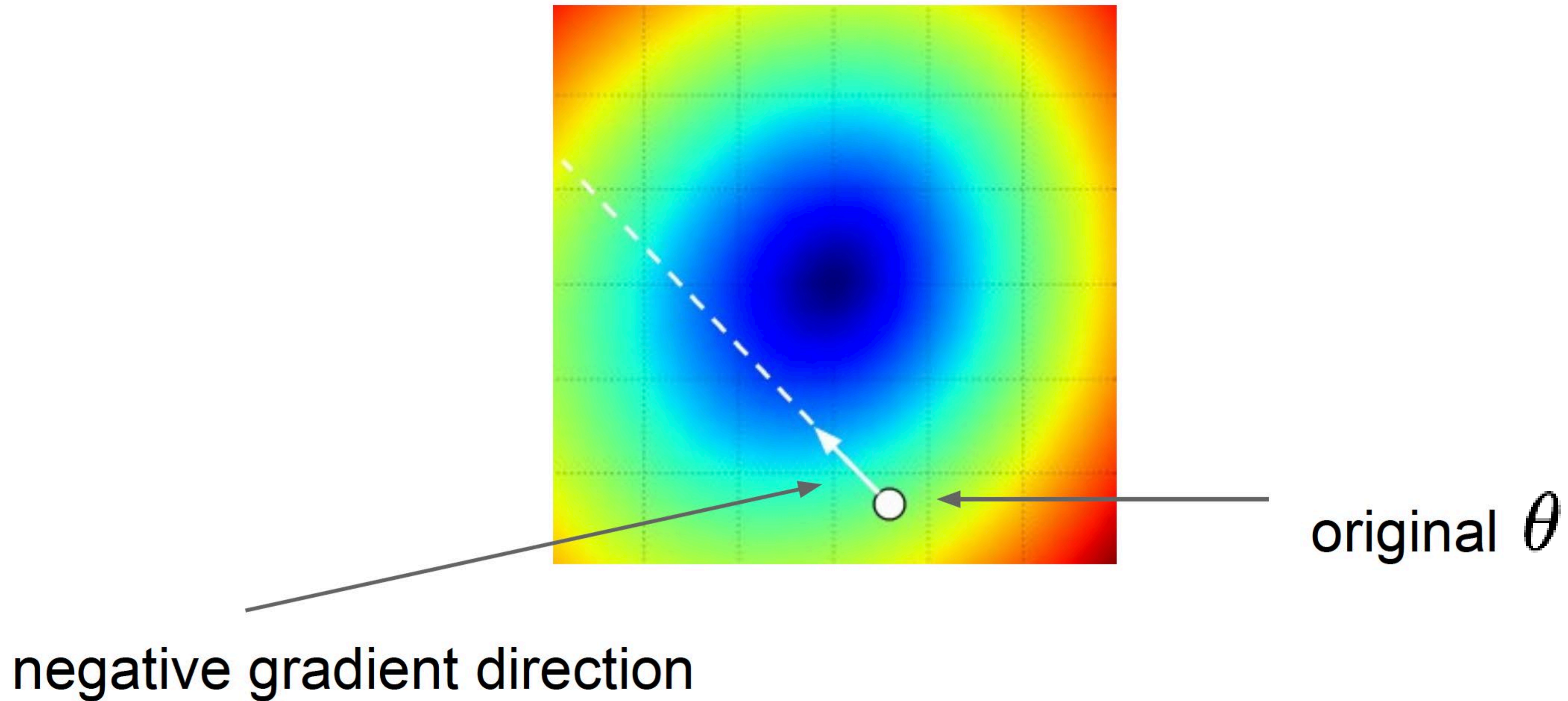
vector of parameter partial derivatives

b. Gradient update

$$\theta \leftarrow \theta - \eta \frac{\partial \mathcal{L}}{\partial \theta}$$

vector of parameter update equations

Learning rates



$$\theta \leftarrow \theta - \eta \frac{\partial \mathcal{L}}{\partial \theta}$$

Step size: learning rate

Too big: will miss the minimum

Too small: slow convergence

Introduction

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad B = \begin{bmatrix} 4 & 3 \\ 2 & 1 \end{bmatrix}$$

$$A + B = \begin{bmatrix} 5 & 5 \\ 5 & 5 \end{bmatrix}$$

$$A - B = \begin{bmatrix} -3 & -1 \\ 1 & 3 \end{bmatrix}$$

$$AB = \begin{bmatrix} 11 & 10 \\ 16 & 14 \end{bmatrix}$$



